

Structural Nested Models and G-estimation: The Partially Realized Promise

Stijn Vansteelandt and Marshall Joffe

Abstract. Structural nested models (SNMs) and the associated method of G-estimation were first proposed by James Robins over two decades ago as approaches to modeling and estimating the joint effects of a sequence of treatments or exposures. The models and estimation methods have since been extended to dealing with a broader series of problems, and have considerable advantages over the other methods developed for estimating such joint effects. Despite these advantages, the application of these methods in applied research has been relatively infrequent; we view this as unfortunate. To remedy this, we provide an overview of the models and estimation methods as developed, primarily by Robins, over the years. We provide insight into their advantages over other methods, and consider some possible reasons for failure of the methods to be more broadly adopted, as well as possible remedies. Finally, we consider several extensions of the standard models and estimation methods.

Key words and phrases: Causal effect, confounding, direct effect, instrumental variable, mediation, time-varying confounding.

1. INTRODUCTION

Structural nested models (SNMs) were designed in part to deal with confounding by variables affected by treatment (Robins, 1986). The problem arises when one is interested in estimating the joint effect of a sequence of treatments in the presence of a variable L with three characteristics, depicted in Figure 1:

1. It is independently associated with the outcome Y of interest. This can happen because (a) it is a direct cause of the outcome, or because (b) it shares unmeasured common causes with the outcome of interest.
2. It predicts subsequent levels (A_1) of the treatment;
3. It is affected by earlier treatment (A_0).

As a motivating example, consider an observational study of the effect of erythropoietin alpha (EPO) on mortality in a population with end-stage renal disease (ESRD) receiving hemodialysis. Patients on dialysis tend to be anemic, as commonly measured via hematocrit (Hct) or hemoglobin levels. EPO is used to treat the anemia and stimulate the body's production of red blood cells; Hct (L) thus satisfies covariate characteristic 3. Furthermore, patients with more severe anemia (lower Hct) typically receive higher doses of EPO (characteristic 2), and sicker patients tend to be more anemic [characteristic 1(b)]. Both these characteristics 1 and 2 make Hct a confounder of the effect of later treatment, requiring adjustment to estimate the effect of EPO A_1 . Observational studies of the effect of extended EPO dosing on mortality will thus be characterized by confounding by a variable (Hct) affected by treatment.

In settings like the above, where the interest lies in estimating the joint effect of a sequence of treatments, standard methods which attempt to estimate these effects simultaneously (e.g., regression of Y on A_0 and A_1 or some function of both) will be inappropriate,

Stijn Vansteelandt is Professor of Statistics, Department of Applied Mathematics, Computer Science and Statistics, Ghent University, B-9000 Gent, Belgium (e-mail: stijn.vansteelandt@UGent.be). Marshall Joffe is Professor of Biostatistics, University of Pennsylvania, Perelman School of Medicine, Philadelphia, USA (e-mail: mjoffe@mail.med.upenn.edu).

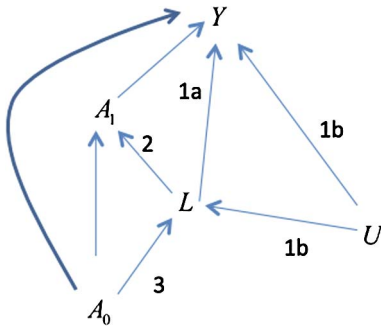


FIG. 1. Causal diagram for time-varying treatment.

whether or not one adjusts for or conditions on the confounder L . Characteristics 1(a) and 3 make Hct (L) an intermediate variable on the pathway from early EPO treatment A_0 to outcome Y ; adjustment for it blocks the path $A_0 \rightarrow L \rightarrow Y$, making it impossible to find the part of the effect of early EPO treatment (A_0) mediated by Hct. Characteristics 1(b) and 3 make Hct (L) a so-called collider (Pearl, 1995) on the path $A_0 \rightarrow L \leftarrow U \rightarrow Y$; conditioning on or adjusting for it induces associations between A_0 and Y even if no effect of A_0 on Y exists.

Over an extended period of time, James Robins (with some help from collaborators) introduced three basic approaches for dealing with such confounding: the parametric G-formula (Robins, 1986), structural nested models (Robins, 1989; Robins et al., 1992) with the associated method of G-estimation and marginal structural models (Robins, Hernan and Brumback, 2000) with the associated method of inverse probability of treatment weighting. As we will argue throughout this paper, SNMs and G-estimation are, in principle, better tailored for dealing with failure of the usual assumptions of no unmeasured confounders or sequential ignorability often used to justify the application of all of these methods, as well as with (near) positivity violations whereby certain strata contain (nearly) no treated or untreated subjects (Robins, 2000). Despite these advantages, the application of these methods in applied research has been relatively infrequent.

Broadly speaking, there are two types of SNMs: models for the effect of a treatment or sequence of treatments on the mean of an outcome, and models for the effect of a treatment on the entire distribution of the outcome(s). The former include structural nested mean models (SNMMs), which have close links to structural nested cumulative failure time models (SNCFTMs) for survival outcomes; the latter include structural nested distribution models (SNDMs), which have close links to structural nested failure time models (SNFTMs) for

survival outcomes. For pedagogic purposes, we will introduce these models first for point treatments (i.e., treatments which are administered at one specific time point) in Section 2. We then discuss identifying assumptions and the associated G-estimation method in Section 3, and contrast it with alternative estimation methods for the effect of a point treatment in Section 4. These results are extended to time-varying treatments in Sections 5 and 6. We show how to predict the effects of interventions in Section 7, examine extensions to mediation analysis in Section 8 and conclude with a discussion.

2. STRUCTURAL MODELS FOR POINT TREATMENTS

2.1 Structural Mean Models

Let Y^a denote the outcome in a given subject that would be seen were the subject to receive treatment a . This variable is a potential outcome, which we connect to the observed outcome through the consistency assumption that $Y = Y^a$ if the observed treatment $A = a$; otherwise, Y^a is counterfactual. Causal effects can now be defined as comparisons of potential outcomes Y^a and Y^{a^\dagger} for the same individual subject or group of subjects for different treatments a and a^\dagger (Rubin, 1978; Robins, 1986). In particular, letting $a^\dagger = 0$ for notational convenience, average causal effects can be defined in terms of comparisons of average potential outcomes, for example, $E(Y^a | L = l, A = a) - E(Y^0 | L = l, A = a)$ or $E(Y^a | L = l, A = a)/E(Y^0 | L = l, A = a)$.

Structural Mean Models (SMMs) (Robins, 1994, 2000) parameterize average causal effects in subjects receiving level a of treatment as

$$(1) \quad g\{E(Y^a | L = l, A = a)\} - g\{E(Y^0 | L = l, A = a)\} = \gamma^*(l, a; \psi^*),$$

for all l and a . Here, $g(\cdot)$ is a known link function (e.g., the identity, log or logit link), $\gamma^*(l, a; \psi)$ is a known function, smooth in ψ and satisfying $\gamma^*(l, 0; \psi) = 0$ for all l and ψ . Here and throughout, ψ^* is the true unknown finite-dimensional parameter. With $a = 0$ encoding absence of treatment—as we will assume throughout—SMMs thus express the effect of removal of treatment on the outcome mean.

Typically, the parameterization is chosen to be such that $\gamma^*(l, a; 0) = 0$ for all a and l , so that $\psi^* = 0$ encodes the null hypothesis of no treatment effect. For instance, for scalar covariate L one may consider the

additive or linear SMM [which uses the identity link $g(x) = x$]:

$$(2) \quad \begin{aligned} & E(Y^a | L = l, A = a) - E(Y^0 | L = l, A = a) \\ &= (\psi_0^* + \psi_1^* l)a, \end{aligned}$$

for unknown ψ_0^*, ψ_1^* . With A a binary exposure coded as 1 for treatment and 0 for no treatment, ψ_0^* thus encodes the average treatment effect in the treated with covariate value $L = 0$, and ψ_1^* measures how much the average treatment effect in the treated differs between subgroups with a unit difference in L . Likewise, the multiplicative or loglinear SMM uses the log link $g(x) = \log(x)$, for example,

$$\frac{E(Y^a | L = l, A = a)}{E(Y^0 | L = l, A = a)} = \exp\{(\psi_0^* + \psi_1^* l)a\},$$

and the logistic SMM uses the logit link $g(x) = \text{logit}(x)$, for example,

$$\frac{\text{odds}(Y^a = 1 | L = l, A = a)}{\text{odds}(Y^0 = 1 | L = l, A = a)} = \exp\{(\psi_0^* + \psi_1^* l)a\},$$

where $\text{odds}(V = 1 | W) \equiv P(V = 1 | W)/P(V = 0 | W)$ for random variables V and W . If treatment A can take on more than two values, then—without additional assumptions—the function $\gamma^*(l, a; \psi^*)$ cannot be interpreted simply as a dose response function. This is because a dose response would contrast outcomes in the same subset at different levels of a [i.e., contrast $E(Y^a | L = l, A = a)$ with $E(Y^{a'} | L = l, A = a)$ for $a \neq a'$], whereas the functions $\gamma^*(l, a; \psi^*)$ and $\gamma^*(l, a'; \psi^*)$ for $a \neq a'$ contrast causal effects for two different groups (namely, those with $A = a$ versus $A = a'$, but the same $L = l$). We will revisit this subtlety in Section 7.

One can use a SMM to construct a variable $U^*(\psi)$ whose mean value (in a subset of individuals with given covariates and treatment) equals the mean outcome that would have been seen had treatment been removed from that subset. Let

$$U^*(\psi) \equiv Y - \gamma^*(L, A; \psi),$$

if $g(\cdot)$ is the identity link,

$$U^*(\psi) \equiv Y \exp\{-\gamma^*(L, A; \psi)\},$$

if $g(\cdot)$ is the log link and

$$(3) \quad \begin{aligned} & U^*(\psi) \\ & \equiv \text{expit}[\text{logit}\{E(Y | L, A)\} - \gamma^*(L, A; \psi)], \end{aligned}$$

if $g(\cdot)$ is the logit link. Then

$$(4) \quad E\{U^*(\psi^*) | L, A\} = E(Y^0 | L, A).$$

This identity will be central to the estimation methods for ψ^* that we will describe in Section 3. We could have defined $U^*(\psi)$ in general—and in particular for the identity and log link—as $U^*(\psi) \equiv g^{-1}[g\{E(Y | L, A)\} - \gamma^*(L, A; \psi)]$. We have avoided doing this for the identity and log links as it makes the definition of $U^*(\psi)$ dependent on the expectation $E(Y | L, A)$, which can be undesirable when this demands additional modeling. However, this (or some alternative) is unavoidable for the logit link. Special estimation methods will therefore be required for logistic SMMs.

SMMs can also be used to describe the effect of a multivariate point treatment. For instance, for a bivariate treatment $A = (A^{(1)}, A^{(2)})'$, one may use a SMM with $\gamma^*(L, A; \psi) = \psi_1 A^{(1)} + \psi_2 A^{(2)} + \psi_3 A^{(1)} A^{(2)}$ to describe the effect of setting both treatments to zero. When primary interest lies in the interaction (ψ_3) between $A^{(1)}$ and $A^{(2)}$ in their effect on the outcome, then one may instead consider the class of less restrictive Structural Mean Interaction Models (Vansteelandt et al., 2008a; Tchetgen Tchetgen, 2012). To guard against misspecification of the main treatment effects, these further relax the SMM restrictions by merely parameterising the contrast between the effects of $A^{(1)}$ when $A^{(2)}$ is set to some value $a^{(2)}$ versus zero (or of the effects of $A^{(2)}$ when $A^{(1)}$ is set to some value $a^{(1)}$ versus 0):

$$(5) \quad \begin{aligned} & g\{E(Y^{a^{(1)}, a^{(2)}} | A = a, L = l)\} \\ & - g\{E(Y^{0, a^{(2)}} | A = a, L = l)\} \\ & - g\{E(Y^{a^{(1)}, 0} | A = a, L = l)\} \\ & + g\{E(Y^{0, 0} | A = a, L = l)\} \\ & = \gamma^*(l, a^{(1)}, a^{(2)}; \psi^*), \end{aligned}$$

for $a = (a^{(1)}, a^{(2)})'$; here, $\gamma^*(l, a^{(1)}, a^{(2)}; \psi)$ is a known function which encodes the interaction between both treatments, and which must be smooth in ψ and satisfy $\gamma^*(l, 0, a^{(2)}; \psi) = \gamma^*(l, a^{(1)}, 0; \psi) = 0$ for all $l, a^{(1)}, a^{(2)}$ and ψ . For instance, the natural choice $\gamma^*(l, a^{(1)}, a^{(2)}; \psi) = \psi a^{(1)} a^{(2)}$ imposes that the interaction between both exposures is the same at all levels of l .

2.2 Structural Distribution Models

When the outcome mean does not adequately summarize the data or the interest lies more broadly in evaluating treatment effects on the outcome distribution, then Structural Distribution Models (SDMs) can be used instead. These are closely related to SMMs,

but instead map percentiles y of the conditional distribution of Y^a , given $L = l$ and $A = a$, into percentiles $\gamma(y, l, a; \psi^*)$ of the conditional distribution of Y^0 , given $L = l$ and $A = a$. In particular, they postulate that

$$(6) \quad F_{Y^0|L=l, A=a}\{\gamma(y, l, a; \psi^*)\} = F_{Y^a|L=l, A=a}(y),$$

for all l and a . As with SMMs, $\gamma(y, l, a; \psi)$ is a known function, smooth in ψ and satisfying $\gamma(y, l, 0; \psi) = y$ for all y, l . With $a = 0$ encoding the absence of treatment, SDMs thus express the effect of removing treatment on the outcome distribution rather than the outcome mean.

Typically the parameterisation of a SDM is chosen to be such that $\gamma(y, l, a; 0) = y$, so that $\psi^* = 0$ encodes the null hypothesis of no treatment effect. For instance, for scalar covariate L , one could assume that

$$(7) \quad \begin{aligned} F_{Y^0|L=l, A=a}(y - \psi_0^*a - \psi_1^*al) \\ = F_{Y^a|L=l, A=a}(y), \end{aligned}$$

for all l and a . This characterizes a location shift model following which the conditional distribution of Y^0 , given L and A , can be obtained by shifting the conditional distribution of Y , given L and A , by $-\psi_0^*A - \psi_1^*AL$. One can use this to construct a variable

$$U(\psi^*) \equiv \gamma(Y, L, A; \psi^*)$$

whose distribution (in a subset of individuals with given covariates and treatment) is the same as that of the outcome that would have been seen had treatment been removed from that subset, in the sense that

$$(8) \quad F_{Y^0|L, A}(y) = F_{U(\psi^*)|L, A}(y);$$

for example, $U(\psi) = Y - \psi_0^*A - \psi_1^*AL$ in the location shift example. This will be useful for the estimation of ψ^* .

SDMs have a stronger variant called rank preserving SDMs (Robins and Tsiatis, 1991), which postulate that

$$Y^0 = \gamma(Y, L, A; \psi^*).$$

For instance, a stronger variant of the location shift model of the previous paragraph assumes that $Y^0 = Y - \psi_0^*A - \psi_1^*AL$. By making a mapping between the potential outcomes themselves (rather than between distributions), such rank preserving SDMs are easier to understand and communicate. However, they are seldom plausible because they impose that the rankings of two subjects with different outcome values but identical treatment and covariates are preserved after mapping into Y^0 (hence the term “rank preserving”). In

particular, they assume that subjects with identical outcome, treatment and covariate values experience identical treatment effects.

Location shift SDMs like (7) make substantially stronger assumptions than correspondingly parameterized SMMs. The distribution models assume that treatment level a shifts each percentile of the conditional distribution of Y , given $L = l, A = a$ by a value $\gamma^*(l, a; \psi^*)$ constant for all y [i.e., $\gamma(y, l, a; \psi) = y - \gamma^*(l, a; \psi)$], whereas the mean model assumes only a mean shift of $\gamma^*(l, a; \psi)$. When location shift is implausible, it can sometimes be made more plausible by transforming y . For instance, for strictly positive y , one might obtain a location shift SDM by defining $\gamma(y, l, a; \psi) = \exp\{\log(y) - \gamma(l, a; \psi)\}$. There will then be a correspondence between the parameters of the SDM and those of a SMM for $\log(y) - \log\{\gamma(y, l, a; \psi)\}$.

The parameterization and interpretation of SDMs that are not simply shift models can be tricky. This is because, by the nature of the cumulative distribution function, the function $\gamma(y, l, a; \psi)$ must be increasing in y for each l, a and ψ , and it may be difficult to impose that. For instance, the function $\gamma(y, a, l; \psi) = y - a\psi_1 - ya\psi_2$ may appear natural, but is not guaranteed increasing in y . An alternative function which is naturally increasing in y is $\gamma(y, a, l; \psi) = y \exp(-a\psi_2) - a\psi_1$. Here, interpretation is somewhat subtle; while ψ_2 expresses the effect of treatment A on the residual variability of Y , it also has implications for the effect of treatment on the mean of Y , and so ψ_1 cannot be interpreted simply as the effect of treatment on the mean outcome, unless $\psi_2 = 0$.

SDMs lend themselves naturally to the analysis of failure times. For instance, consider model (6) with T^a, T^0 and t substituting for Y^a, Y^0 and y . Then the choice $\gamma(t, a, l; \psi) = t \exp(a\psi_0 + al\psi_1)$ implies the failure time model defined by

$$S_{T^0|L=l, A=a}\{t \exp(-a\psi_0^* - al\psi_1^*)\} = S_{T|L=l, A=a}(t),$$

for all l and a , where $S(\cdot)$ denotes the survival function. This model, which is an example of a Structural Accelerated Failure Time Model (SAFTM) (Robins, 1989; Robins and Tsiatis, 1991; Robins, 1992; Robins et al., 1992), expresses that treatment lengthens lifetime by a factor $\exp(a\psi_0^* + al\psi_1^*)$ (in distribution) among subjects with treatment a and covariate l .

2.3 Structural Mean and Distribution Models for Repeated Measures Outcomes

Structural mean and distribution models require some modification for repeated measures outcomes.

The modifications for SMMs are simpler, but also allow a new class of models for discrete-time failures. Extension of SDMs is more complicated. We consider these in order.

We begin with some notation common to both types of models. Suppose that measurements on exposure and confounders are collected at time point t_0 and that outcome measurements are recorded at fixed later time points t_1, \dots, t_{K+1} . Let for a variable X , X_k denote the level of the variable that one obtains at time t_k . We use overbars to denote the history of a variable; thus, $\bar{X}_k = \{X_0, X_1, \dots, X_k\}$ denotes the history of X through t_k . We use underbars to denote the future of a variable; thus, $\underline{X}_k \equiv \{X_k, \dots, X_{K+1}\}$. Finally, we use \underline{X} as shorthand notation for \underline{X}_1 and $X_{k:m}$ for $m \geq k$ to denote (X_k, \dots, X_m) .

2.3.1 Structural mean models and structural cumulative failure-time models. Extension of SMMs to repeated measures is relatively straightforward, because they model separately the effect of a treatment on each component outcome. SMMs parameterize contrasts of \underline{Y}^a and \underline{Y}^0 as

$$g\{E(\underline{Y}^a | L = l, A = a)\} - g\{E(\underline{Y}^0 | L = l, A = a)\} = \gamma^*(l, a; \psi^*),$$

for all l and a . Here, $g(\cdot)$ is a known $(K + 1)$ -dimensional link function, $\gamma^*(l, a; \psi)$ is a known $(K + 1)$ -dimensional function with components $\gamma_k^*(l, a; \psi)$, $k = 1, \dots, K + 1$, that parameterize the effect of treatment on Y_k . These components are smooth in ψ and satisfy $\gamma_k^*(l, 0; \psi) = 0$ for all l and ψ . For instance, the SMM defined by

$$E(Y_k^a | L = l, A = a) - E(Y_k^0 | L = l, A = a) = (\psi_0^* + \psi_1^* l) a (t_k - t_0),$$

for $k = 1, \dots, K + 1$, expresses that the effect of treatment a may depend on covariates l and changes linearly over time, being zero at the baseline time t_0 .

Under this repeated measures SMM, as in Section 2.1, it is possible to define a transformation $U^*(\psi)$ of the observed outcome vector \underline{Y} so that

$$E\{U^*(\psi^*) | L, A\} = E(\underline{Y}^0 | L, A).$$

Here, $U^*(\psi)$ is a vector with components $Y_k - \gamma_k^*(L, A; \psi)$ for $k = 1, \dots, K + 1$ if $g(\cdot)$ is the identity link, $Y_k \exp\{-\gamma_k^*(L, A; \psi)\}$ if $g(\cdot)$ is the log link, and $\text{expit}[\text{logit}\{E(Y_k | L, A)\} - \gamma_k^*(L, A; \psi)]$ if $g(\cdot)$ is the logit link.

Structural Cumulative Failure Time Models (SCFTMs; Picciotto et al., 2012) are a variant of repeated measures loglinear SMMs for the modeling of cumulative failure time probabilities:

$$\frac{P(T^a < t_k | L = l, A = a)}{P(T^0 < t_k | L = l, A = a)} = \exp\{\gamma_k^*(l, a; \psi^*)\},$$

for all l, a and $k = 1, \dots, K + 1$. A limitation of this class of models is that their parameterization can be tricky when the cumulative probability of failure becomes large, because the model does not restrict the outcome probabilities to stay below 1. Martinussen et al. (2011) independently proposed a continuous-time version of this model and lay out connections with additive hazard models.

2.3.2 Structural distribution models. For multivariate outcomes, SDMs parameterize the effect of a treatment A on the marginal distribution of the vector of future potential outcomes \underline{Y}^a . This mapping is typically done recursively, taking the components Y_k^a and Y_k in forward sequence. These models are therefore most easily understood by first considering the class of more restrictive rank-preserving SDMs, which postulate that, for subjects with $A = a$ and $L = l$:

$$(9) \quad Y_k^0 = \gamma_k(Y_k, \bar{Y}_{k-1}, l, a; \psi^*)$$

for $k = 1, \dots, K + 1$. Here, $\gamma_k(y_k, \bar{y}_{k-1}, l, a; \psi)$ is a known function, smooth in ψ and monotonic in y_k , and $\gamma_k(y_k, \bar{y}_{k-1}, l, 0; \psi) = y_k$ for all \bar{y}_{k-1}, l , and ψ . For instance, with two time points ($K = 1$), a rank preserving SDM may be given by the following set of restrictions:

$$(10) \quad \begin{aligned} Y_2^0 &= Y_2 - (\psi_1^* + \psi_2^* Y_1) A, \\ Y_1^0 &= Y_1 - \psi_3^* A. \end{aligned}$$

If the effect of A on Y_2 varies with Y_1 , as in this example, then one must model this explicitly since the model would otherwise—perhaps unrealistically—assume that treatment does not affect the correlation between repeated outcomes (conditional on A, L). This is unlike in SMMs where one can average the effect of A on Y_2 over all Y_1 -values. This makes it substantially more difficult to parameterize SDMs than SMMs. It moreover complicates the interpretation of effects; for example, ψ_1^* in (10) is difficult to interpret when $\psi_2^* \neq 0$ since it expresses the effect of treatment on Y_2 in subjects with $A = 1$ and $Y_1 = 0$, where Y_1 may itself be affected by treatment. Equation (10) may hence be easier to interpret upon re-expressing it as

$$\begin{aligned} Y_2^0 &= Y_2 - \{\psi_1^* + \psi_2^* (Y_1^0 + \psi_3^* A)\} A \\ &= Y_2 - (\psi_1^* + \psi_2^* Y_1^0 + \psi_2^* \psi_3^* A) A. \end{aligned}$$

A SDM relaxes the restrictions of the rank-preserving SDM by demanding that the equality (9) merely holds in distribution, conditional on $L = l$ and $A = a$. Assuming that given L and A , \underline{Y} has a continuous multivariate distribution with probability 1, a SDM can thus be defined by the set of restrictions

$$\begin{aligned} F_{\underline{Y}^0|L=l, A=a} \{ \gamma(\underline{y}, l, a; \psi^*) \} &= F_{\underline{Y}^a|L=l, A=a}(\underline{y}) \\ &= F_{\underline{Y}|L=l, A=a}(\underline{y}), \end{aligned}$$

for all l, a , where

$$\begin{aligned} \gamma(\underline{y}, l, a; \psi^*) &\equiv \{ \gamma_1(y_1, l, a; \psi^*), \\ &\quad \gamma_2(\bar{y}_2, l, a; \psi^*), \dots, \\ &\quad \gamma_{K+1}(\bar{y}_{K+1}, l, a; \psi^*) \} \end{aligned}$$

is given by (Robins, Rotnitzky and Scharfstein, 2000):

$$\begin{aligned} \gamma_1(y_1, l, a; \psi^*) &= F_{Y_1^0|L=l, A=a}^{-1} \circ F_{Y_1|L=l, A=a}(y_1), \\ \gamma_k(\bar{y}_k, l, a; \psi^*) &= F_{Y_k^0|L=l, A=a, \bar{Y}_{k-1}^0 = \gamma_{1:k-1}(\bar{y}_{k-1}, l, a; \psi^*)}^{-1} \\ &\quad \circ F_{Y_k|L=l, A=a, \bar{Y}_{k-1} = \bar{y}_{k-1}}(y_k), \end{aligned}$$

for $k = 2, \dots, K + 1$. For instance, the SDM corresponding to (10) may be written:

$$\begin{aligned} &F_{Y_1^0|L=l, A=a}(y_1 - \psi_3^* a) \\ &= F_{Y_1|L=l, A=a}(y_1), \\ (11) \quad &F_{Y_2^0|L=l, A=a, Y_1^0 = y_1 - \psi_3^* a} \{ y_2 - (\psi_1^* + \psi_2^* y_1) a \} \\ &= F_{Y_2|L=l, A=a, Y_1 = y_1}(y_2). \end{aligned}$$

The decomposition of the causal effects in the blip functions $\gamma_k(\bar{y}_k, l, a; \psi^*)$ is recursive because one must model not merely average effects but instead the full mapping between distributions. In particular, the effect of treatment on the first potential outcome is modeled first; then, mapping between distributions is done successively for the outcome at successive times. The overall blip function encoded by $\gamma(\cdot)$ and the first element of this function has the usual structure and interpretation of causal estimands; that is, as a comparison of distributions of potential outcomes under different interventions for the same group of subjects. However, the component functions $\gamma_k(\cdot)$, $k > 1$ do not in general have this interpretation, since the conditioning in these mapping functions is not common between Y_k^0 and Y_k ; for instance, the left-hand side of (11) conditions on Y_1^0 , whereas the right-hand side conditions on Y_1 . Nonetheless, these component functions are causal in the sense that they represent the impact of treatment on the conditional distribution of a variable. This

feature is shared with the causal rate or hazard ratio (Hernan, 2010). Under the strong assumption of rank preservation, the conditioning is on a common variable, and so then the components of the blip function do have a standard causal interpretation.

For repeated measures outcomes, SDMs correspond with similarly parameterized SMMs if the SDMs are shift models. In a shift SDM, the component functions $\gamma_k(y_k, \bar{y}_{k-1}, l, a)$ may be written as $\gamma_k(y_k, \bar{y}_{k-1}, l, a) = y_k - \gamma_k^*(l, a; \psi)$. These require that the shift in percentiles of the distribution of y_k not only be independent of y_k but also of \bar{y}_{k-1} . Thus, shift SDMs make substantially stronger assumptions than similarly parameterized SMMs.

Under the SDM, a $(K + 1)$ -dimensional variable $U(\psi^*) = \{U_1(\psi^*), \dots, U_{K+1}(\psi^*)\}$ can be constructed with components $U_k(\psi) = \gamma_k(\bar{Y}_k, L, A; \psi)$. This vector mimics the counterfactual outcome vector \underline{Y}^0 in the sense that

$$P\{U(\psi^*) > \underline{y} \mid L, A\} = P(\underline{Y}^0 > \underline{y} \mid L, A).$$

This result will be useful for estimation.

2.4 Retrospective Blip Models

The blip functions and causal models discussed above largely consider the effect of a blip of treatment conditional only on treatment and baseline covariates; the sole exception has been SDMs for repeated measures outcomes, where the effect of treatment on later outcomes is modeled additionally conditional on earlier outcomes, and where the interpretation of the model parameters is not clear as a usual causal contrast. This focus is consistent with an orientation of the models to be more directly useful for making decisions, where the effect of treatment is modeled conditional only on information available at the time of the decision.

For explanatory purposes, one can construct structural models for the effect of a treatment conditional on information not available at the time of treatment. Such models may have explanatory uses even though the quantities they model are less directly relevant for making decisions. Consider modeling the effect of screening mammography on breast cancer mortality (Joffe, Small and Hsu, 2007). To a first approximation, one might assume that the mammogram has an effect on death only among subjects for whom it detects a tumor. Suppose that some subjects undergo screening at the start of the study ($A = 1$; $A = 0$ otherwise). Let L_1 indicate 1 if cancer is detected at time t_1 after the start of the study and 0 otherwise. It is of interest to know

how much the screening mammogram affects mortality for subjects for whom it is effective in detecting cancer. We can then model the effect of the treatment on the outcome using a retrospective SDM (RSDM) or SFTM, which conditions on L_1 in addition to treatment and baseline covariates:

$$(12) \quad \begin{aligned} F_{Y^0|L_0=l_0, L_1=l_1, A=a} \{ \gamma(y, l_0, l_1, a; \psi^*) \} \\ = F_{Y|L_0=l_0, L_1=l_1, A=a}(y). \end{aligned}$$

In this example, we might assume that $\gamma(y, l_0, 0, a; \psi^*) = y$ to reflect that screening has no effect in subjects for whom no tumor is detected. Note that though L_1 may be affected by A , conditioning on it does not distort the interpretation of the parameters as encoding a causal effect because identity (12) still involves a comparison of the same subjects (those with $L_0 = l_0, L_1 = l_1, A = a$) under different interventions.

Models of this sort might also be useful in determining whether the effect of a treatment given at baseline is modified by post-treatment covariates and so whether there are identifiable subgroups of subjects for whom treatment appears not to be working (Stephens, Keele and Joffe, 2013). Changes or additions to treatment might then be proposed in such subgroups after baseline. Joffe, Small and Hsu (2007) consider the relation between these retrospective models and the popular approach of principal stratification (Frangakis and Rubin, 2002). These models can generalize to a sequence of time-varying treatments, where there are additional justifications for their use (see Section 5.3).

3. IDENTIFICATION AND ESTIMATION IN STRUCTURAL MODELS FOR POINT TREATMENTS

Two kinds of assumptions have been proposed for use in most of the literature on estimation in SMMs and SDMs: no unmeasured confounders and instrumental variables type assumptions. In this section, we will focus on the former, and defer discussion of the latter to Section 6.3.

3.1 Ignorability

The required no unmeasured confounders assumption for the identification of the parameter ψ^* indexing SMMs and SDMs can be formulated as

$$(13) \quad A \perp\!\!\!\perp Y^0 \mid L,$$

where $U \perp\!\!\!\perp V \mid W$ for random variables U, V, W denotes that U is conditionally independent of V , given W . This assumption, which is empirically unverifiable,

expresses that L is sufficient to adjust for confounding of the association between A and Y . Assumption (13), which is also referred to as the weak ignorability or exchangeability assumption, is weaker than the strong ignorability assumption of Rosenbaum and Rubin (1984) which, for binary treatments, states that $A \perp\!\!\!\perp (Y^0, Y^1) \mid L$. However, it is generally difficult to imagine settings where assumption (13) holds, but strong ignorability fails (one exception might be settings where individuals choose treatment on the basis of their perceived belief of benefit, which may be correlated with actual benefit $Y^1 - Y^0$). That (13) is a weaker assumption is exhibited in the fact that, for binary treatments, it only identifies the effect of treatment on the treated—a contrast that has been of interest in econometrics and epidemiology (Greenland and Robins, 1986):

$$\begin{aligned} E(Y^1 - Y^0 \mid A = 1, L) \\ = E(Y^1 \mid A = 1, L) - E(Y^0 \mid A = 1, L) \\ = E(Y^1 \mid A = 1, L) - E(Y^0 \mid A = 0, L) \\ = E(Y \mid A = 1, L) - E(Y \mid A = 0, L); \end{aligned}$$

the second equality follows due to ignorability (13) and the third due to the consistency assumption. The parameters of SMMs, SDMs and SCFTMs represent the effect of treatment in the treated (or, more generally, the effect of receiving treatment level a for subjects who received level a of treatment), and so this weaker assumption is sufficient for identification.

It follows by a similar reasoning that the blip functions in the SMMs and SDMs discussed in Sections 2.1–2.3 are nonparametrically just identified under ignorability (Robins, Rotnitzky and Scharfstein, 2000). That is, the contrast of the outcomes under the observed treatment and the outcomes that would have been seen in the absence of treatment is computable for each level of a and l (and, for SDMs, of y) from the law of the observables without assuming any restrictions or parameterization on these functions. While such nonparametric identification is of limited use in complex settings (especially with time-varying treatments considered subsequently), due to the curse of dimensionality (Robins and Ritov, 1997), it does ensure the ability to check the assumptions in any assumed causal model (provided a sufficient sample size). In contrast, the retrospective blip functions considered in Section 2.4 are not identified nonparametrically (Vansteelandt, 2010; Stephens, Keele and Joffe, 2013). Multiple retrospective blip models may thus explain the same law of the observables equally well even under ignorability.

3.2 Estimation Under Ignorability

The SMM together with the ignorability assumption (13) implies that

$$\begin{aligned} E\{U^*(\psi^*) | L, A\} &= E(Y^0 | A, L) = E(Y^0 | L) \\ &= E\{U^*(\psi^*) | L\}. \end{aligned}$$

Estimation of ψ^* in a SMM can thus be based on solving estimating equations:

$$(14) \quad 0 = \sum_{i=1}^n [d^*(A_i, L_i) - E\{d^*(A_i, L_i) | L_i\}] \cdot [U_i^*(\psi) - E\{U_i^*(\psi) | L_i\}],$$

which essentially set the empirical conditional covariance between $U^*(\psi)$ and arbitrary functions $d^*(A, L)$ of the dimension of ψ , given L , to zero. For instance, for model (2), the choice $d^*(A_i, L_i) = (1, L_i)'A_i$ results in estimating equations

$$(15) \quad 0 = \sum_{i=1}^n \left(\frac{1}{L_i} \right) \{A_i - E(A_i | L_i)\} \cdot [Y_i - E(Y_i | L_i) - (\psi_0 + \psi_1 L_i)\{A_i - E(A_i | L_i)\}],$$

from which estimates for (ψ_0, ψ_1) can be solved. A locally efficient estimator of ψ^* [under the SMM together with the ignorability assumption (13)] can be attained by setting

$$d^*(A, L) = E \left\{ \frac{\partial U^*(\psi^*)}{\partial \psi} \mid A, L \right\},$$

when the variance of $U^*(\psi^*)$ given A, L is constant; local here means that the efficiency is only attained when this constant variance assumption is met and models for all conditional expectations involved in (14) are correctly specified.

The SDM together with the ignorability assumption (13) implies the more restrictive constraint that

$$(16) \quad U(\psi^*) \perp\!\!\!\perp A \mid L.$$

This motivates estimating ψ^* by picking the value ψ that makes this conditional independence hold. This forms the default approach in SAFTMs, where estimation is based on a grid search whereby the independence (16) is tested for different values of ψ^* using a (standard) statistical test until it is found to be satisfied (Robins et al., 1992). Equivalently, estimation can be

based on solving an estimating equation of the form

$$(17) \quad 0 = \sum_{i=1}^n d\{U_i(\psi), A_i, L_i\} - E[d\{U_i(\psi), A_i, L_i\} | L_i, U_i(\psi)] - E(d\{U_i(\psi), A_i, L_i\} - E[d\{U_i(\psi), A_i, L_i\} | L_i, U_i(\psi)] | A_i, L_i),$$

for ψ , where $d\{U_i(\psi), A_i, L_i\}$ is an arbitrary index function of the dimension of ψ ; for example, $d\{U_i(\psi), A_i, L_i\} = (1, L_i)'A_i U_i(\psi)$. A locally efficient estimator of ψ^* [under the SDM together with the ignorability assumption (13)] can be obtained by solving (17) with $d\{U(\psi), A, L\} = E\{S_\psi(\psi) | U(\psi), A, L\}$, where $S_\psi(\psi)$ is the score for ψ under the observed data likelihood

$$(18) \quad \frac{\partial U(\psi^*)}{\partial Y} f(L) f\{U(\psi^*) | L\} f(A | L)$$

with all components substituted by suitable parametric models (Robins, 1997). For instance, under model (7) with $U(\psi^*)$ given L following a normal distribution with mean linear in L and constant variance, $S_\psi(\psi) = (1, L)'A\{aU(\psi) + bL + c\}$ for certain constants a, b, c , so that a locally efficient estimator is obtained by solving (15).

Estimating equations of form (14) and (17) may also be used for repeated measures outcomes. In (14), $d^*(A_i, L_i)$ now becomes a $p \times (K + 1)$ -dimensional matrix, with p the dimension of ψ . In (17), $d\{U_i(\psi), A_i, L_i\}$ remains an arbitrary index function of the dimension of ψ ; for example, $d\{U_i(\psi), A_i, L_i\} = (1, L_i)'A_i \sum_{m=1}^{K+1} U_{im}(\psi)$.

REMARK. Note that the SMM together with assumption (13) is the same model for the observables as the semiparametric regression model (Chamberlain, 1987):

$$(19) \quad g\{E(Y | L, A)\} = \omega(L) + \gamma^*(L, A; \psi^*),$$

with $\omega(L)$ unspecified. Likewise, the SAFTM [with, e.g., $\gamma(t, a, l; \psi) = t \exp(-a\psi)$] together with assumption (13) can be viewed as a semiparametric generalization of the accelerated failure time model (Wei, 1992), defined by $\log T = \psi A + \varepsilon$ with $\varepsilon \perp\!\!\!\perp A \mid L$.

Because of the curse of dimensionality, evaluating the conditional expectations appearing in equations (14) and (17) requires a parametric working model \mathcal{A} for the conditional distribution of the exposure A :

$$f(A | L) = f(A | L; \alpha^*);$$

here $f(A | L; \alpha)$ is a known density function, smooth in α , and α^* is an unknown finite-dimensional parameter. For instance, for dichotomous exposure, one could assume that $P(A = 1 | L) = \text{expit}(\alpha_0^* + \alpha_1^* L)$ with $\alpha^* = (\alpha_0^*, \alpha_1^*)'$. Here, α^* can be estimated via standard (maximum likelihood) methods.

Evaluating (14) and (17) moreover requires a parametric working model \mathcal{B} for the conditional distribution of $U(\psi^*)$ or the conditional expectation of $U^*(\psi^*)$. For (17), we model:

$$f\{U(\psi^*) | L\} = f\{U(\psi^*) | L; \beta^*\},$$

where $f\{U(\psi^*) | L; \beta\}$ is a known density function, smooth in β , and β^* is an unknown finite-dimensional parameter; to evaluate equations (14), specification of the conditional mean of $U^*(\psi^*)$, given L , suffices. For instance, for a continuous outcome, one could assume that conditional on L and for given ψ^* , $U(\psi^*) = Y - \psi_0^* A - \psi_1^* AL$ is normally distributed with mean $\beta_0^* + \beta_1^* L$ and variance β_2^{*2} , with $\beta^* = (\beta_0^*, \beta_1^*, \beta_2^*)'$. For each fixed value of ψ^* , β^* can be estimated using standard regression methods.

A consistent estimator of ψ^* indexing the SMM or SDM can now be obtained by solving equations (14) or (17), respectively, with α^* and β^* substituted by consistent estimators under models \mathcal{A} and \mathcal{B} , respectively. The resulting estimator of ψ^* is called a G-estimator. In SDMs and linear or loglinear SMMs, it has the attractive property of being doubly robust (Robins and Rotnitzky, 2001): consistent when either model \mathcal{A} or model \mathcal{B} is correctly specified (in addition to a correctly specified structural model and ignorability); it does not require both to be correctly specified, nor does it require specifying which of both is correctly specified. That the solution to equation (14) is doubly robust can be seen because this equation has mean zero at $\psi = \psi^*$ when either model \mathcal{A} or model \mathcal{B} is correctly specified, even if one of them is misspecified. Equation (17) is likewise seen to have mean zero at $\psi = \psi^*$ under model \mathcal{B} ; that it also has mean zero under model \mathcal{A} at $\psi = \psi^*$ is seen by rewriting the equation as

$$\begin{aligned} 0 = & \sum_{i=1}^n d\{U_i(\psi), A_i, L_i\} \\ & - E[d\{U_i(\psi), A_i, L_i\} | L_i, A_i] \\ & - E(d\{U_i(\psi), A_i, L_i\} \\ & - E[d\{U_i(\psi), A_i, L_i\} | L_i, A_i] | U_i(\psi), L_i). \end{aligned}$$

The result now follows, provided that the parameters α and β are variation-independent (i.e., not functionally related), so that a consistent estimator of α^* does

not require consistent estimation of β^* and vice versa. Sandwich standard errors are obtained via the usual estimating equations theory.

In logistic SMMs, to the best of our knowledge, no estimators of ψ^* have been found that are root- n consistent under model \mathcal{A} and the ignorability assumption. This is because the evaluation of $U^*(\psi)$ is anyway dependent upon a model for the conditional mean $E(Y | A, L)$ [see (3)]. Tchetgen Tchetgen, Robins and Rotnitzky (2010) show that double robustness can instead be attained against misspecification of either a model for the density $f(Y | A = 0, L)$ or a model for the density $f(A | Y = 0, L)$. Their key to estimation of ψ^* is that the parameterized association $\gamma^*(L, A; \psi)$ between A and Y , when evaluated at $\psi = \psi^*$, can be used to render A and Y conditionally independent (given L) via inverse probability weighting. Their results apply equally to case-control designs (Tchetgen Tchetgen and Rotnitzky, 2011).

For Structural Mean Interaction Models, inference is developed in Vansteelandt et al. (2008a) when $g(\cdot)$ is the identity or log link and in Tchetgen Tchetgen (2012) when $g(\cdot)$ is the logistic link. Tchetgen Tchetgen and Robins (2010) focus on case-only designs and note that when $g(\cdot)$ is the log link, the multiplicative interaction (5) is identical to the conditional odds ratio between $A^{(1)}$ and $A^{(2)}$, given L within the subgroup of cases. This enables the use of results on logistic SMMs (Tchetgen Tchetgen, Robins and Rotnitzky, 2010) for robust estimation of multiplicative interactions under outcome-dependent sampling.

3.3 Censoring

Censoring presents additional challenges for the analysis of failure-time outcomes T . Random censoring or loss to follow-up can be dealt with through inverse probability of censoring weighting (Robins et al., 1992). Type I censoring, also known as censoring by end of follow-up, can be ignored in the analysis of SCFTMs, but must be dealt with in a different fashion in the analysis of SAFTMs. This is because $U(\psi^*)$ involves the failure-time itself, which is missing for all subjects who fail after planned end-of-follow-up; the coarsening process is informative here as it depends on the actual failure time. We will next describe how Type I censoring can be dealt with in the analysis of SAFTMs.

Let C denote the planned end of follow-up time for given individual. C is known for all subjects, even those observed to fail. However, $U(\psi)$ cannot be evaluated for those who do not fail prior to time C . Knowing that $U(\psi^*) \perp\!\!\!\perp A | L$ under ignorability, the aim is

then to find a function $q\{U(\psi), C\}$ which is observable for all individuals and for which

$$q\{U(\psi^*), C\} \perp\!\!\!\perp A \mid L.$$

If such function is found, then ψ^* can be estimated by solving the original estimating equations for SDMs with $q\{U(\psi), C\}$ replacing $U(\psi)$. A natural choice would be $q\{U(\psi), C\} = \min\{U(\psi), U(C, A, L; \psi)\}$ with $U(C, A, L; \psi)$ the blipped-down censoring time, which is defined like $U(\psi)$ but with T substituted by C . However, this choice would not satisfy the required conditional independence property. The reason is that since C is fixed by design, $U(C, A, L; \psi)$ will in general be a function of A when $\psi \neq 0$ and so will generally fail to be conditionally independent of A , given L . Robins and Tsiatis (1991) thus propose to eliminate the dependence of $U(C, A, L; \psi)$ on A by redefining it to be $C(\psi) \equiv \min_a\{U(C, a, L; \psi)\}$. By thus minimizing over all feasible treatments a , any dependence on the observed treatment is broken so that $X(\psi) \equiv \min\{U(\psi), C(\psi)\}$ and $\Delta(\psi) \equiv I\{U(\psi) < C(\psi)\}$ become always observable quantities that are independent of A given L under ignorability, when evaluated at ψ^* . We may thus choose $q\{U(\psi), C\}$ to be an arbitrary function of $X(\psi)$ and $\Delta(\psi)$.

With each choice of $q\{U(\psi), C\}$, some subjects who are observed to fail may be treated as censored when $\psi \neq 0$. This can happen because for some subjects, $C(\psi)$ may be smaller than $U(\psi)$ even though $T < C$. Such subjects are called artificially censored. Artificial censoring has several consequences. Besides decreasing information about ψ^* as more subjects are artificially censored, the estimating equations are not, in general, continuous in ψ . This is because the functions $q\{U(\psi), C\}$ are not generally continuous in ψ , which happens in part because $\Delta(\psi)$ is not a smooth function of ψ . This can present problems for optimization, especially when ψ is a vector, and may moreover imply that the estimating equations have no solution in finite samples. This problem may be mitigated by choosing $q\{U(\psi), C\}$ to be a smooth function of ψ , for example, $q\{U(\psi), C\} = \Delta(\psi)w_\alpha\{X(\psi)/C(\psi)\}$, where $w_\alpha(t) \equiv I(t > 1 - \alpha)(1 - t)/\alpha + I(t \leq 1 - \alpha)$ (Joffe, Yang and Feldman, 2012). Vock et al. (2013) consider functions $q(\cdot; \psi)$ whose first derivatives exist for all ψ ; they appear to have had better success in convergence for their optimization algorithm.

4. PROPERTIES OF G-ESTIMATION IN STRUCTURAL MODELS FOR POINT TREATMENTS UNDER IGNORABILITY

4.1 Comparison with Ordinary Regression Estimators

Insight into the behavior of G-estimators can be garnered by focusing on the simple model \mathcal{M}_{SMM} defined by the ignorability assumption that $Y^a \perp\!\!\!\perp A \mid L$ for $a = 0, 1$, known treatment mechanism $f(A \mid L)$ and the SMM

$$E(Y^a - Y^0 \mid A = a, L) = \psi^*a.$$

Under homoscedasticity (i.e., when the conditional variance of the outcome, given A and L , is a constant σ^2), the locally efficient G-estimator of ψ^* under model \mathcal{M}_{SMM} has influence function (Newey, 1990)

$$(20) \quad E\{\text{Var}(A \mid L)\}^{-1}\{A - E(A \mid L)\} \cdot \{Y - \psi^*A - E(Y - \psi^*A \mid L)\};$$

it can thus in particular be obtained by setting the sample average of these influence functions to zero and solving for ψ^* . For binary treatment A , linear regression adjustment for the propensity score (Rosenbaum and Rubin, 1984) results in an estimator of ψ^* with influence function of the same form (20), but with $E(Y - \psi^*A \mid L)$ substituted by the population least squares fit from a regression of $Y - \psi^*A$ on the propensity score $E(A \mid L)$. Linear regression adjustment for the propensity score can therefore be viewed as an inefficient and nondoubly robust G-estimation approach (Robins, Mark and Newey, 1992). The close relation between G-estimation and regression adjustment for the propensity score is not maintained in nonlinear models, where propensity score adjustment may not only demand correct models for the propensity score, but also for its association with outcome (Vansteelandt and Daniel, 2014). In nonlinear models, due to non-collapsibility of the treatment effect parameter (Greenland, Robins and Pearl, 1999), its meaning may also change depending on whether covariates are adjusted for in addition to the propensity score.

Ordinary regression estimators [in particular, maximum likelihood estimators obtained by fitting model (19) under a finite-dimensional parameterization of $\omega(L)$] are at least as efficient as the previously considered G-estimators, provided correct model specification. From the variance of the influence functions,

we can deduce that the asymptotic variance of the locally efficient G-estimator is

$$(21) \quad \frac{\sigma^2}{E\{\text{Var}(A | L)\}},$$

when there is homoscedasticity and the conditional mean $E(Y - \psi^* A | L) = E(Y | A = 0, L)$ is correctly specified. The ordinary least squares (OLS) estimator under the linear regression model $E(Y | A, L) = \beta' L + \psi A$ has an asymptotic variance which is smaller but, interestingly, usually not much smaller:

$$\frac{\sigma^2}{E[\text{Var}(A | L) + \{E(A | L) - \tilde{E}(A | L)\}^2]}.$$

This follows from its influence function, which is of the same form (20), but with $E(A | L)$ substituted by $\tilde{E}(A | L)$, the population least squares fit from a regression of A on L .

Despite their greater efficiency, ordinary regression estimators have a number of limitations not shared by G-estimators, an important one being their lack of extensibility to the analysis of sequential treatments (see Section 5). Furthermore, their explicit reliance on a model for the association between outcome and covariates can be disadvantageous when the treated and untreated subjects are very different in their covariate distributions, for then even well-fitting models for the outcome may be prone to extrapolation bias (Rosenbaum and Rubin, 1984). This is not the case for G-estimators when they are based on a correctly specified model (\mathcal{A}) for the treatment process. This is also seen from the form of the influence functions (20), following which individuals in regions of little or no overlap [i.e., at covariate values L where $\text{Var}(A | L)$ is small] will hardly contribute in the calculation of the G-estimator because $A - E(A | L) \approx 0$ for such individuals. As with other estimation approaches based on propensity score adjustment (e.g., matching), the information about ψ^* will thus come primarily from regions with sufficient overlap, which we view as desirable. In contrast, OLS estimators are more susceptible to extrapolation bias since the leading term $A - \tilde{E}(A | L)$ in their influence functions may be far from zero for individuals in regions of little or no overlap. Finally, an advantage of G-estimation methods is that they can incorporate a priori knowledge on the exposure distribution. For instance, Vansteelandt et al. (2008b) exploit knowledge on the distribution of offspring genotypes given parental genotypes (based on Mendel's law of segregation), by using G-estimators to develop gene-environment interaction tests that are robust against misspecification of the effect of environmental exposures on the outcome.

4.2 Comparison with Inverse Probability Weighted Estimators

For the analysis of sequential treatments (see Section 5), marginal structural models (MSM) (Robins, Hernan and Brumback, 2000) and inverse probability weighted (IPW) estimators are much more popular than SMMs and SDMs and G-estimators. This is related to G-estimation being computationally more demanding by the lack of off-the-shelf software. It is thus of interest to compare the behaviour of these estimators in a simple setting with dichotomous treatment. Consider therefore model \mathcal{M}_{MSM} , which is defined by the ignorability assumption that $Y^a \perp\!\!\!\perp A | L$ for $a = 0, 1$, known propensity score $E(A | L)$ and the nonparametric MSM

$$E(Y^a) = \alpha + \psi^* a.$$

Note, since $Y^a \perp\!\!\!\perp A | L$ for $a = 0, 1$, that $\psi^* = E(Y^1 - Y^0)$ in both models \mathcal{M}_{SMM} and \mathcal{M}_{MSM} , and thus defines the same parameter. Nonetheless, model \mathcal{M}_{MSM} is less restrictive than model \mathcal{M}_{SMM} in that it does not postulate that the treatment effect is homogeneous (i.e., constant over levels of L). This explains why the asymptotic variance of the locally efficient IPW estimator under model \mathcal{M}_{MSM} , which has influence function (Robins, Rotnitzky and Zhao, 1994)

$$\begin{aligned} & \frac{A\{Y - E(Y | A = 1, L)\}}{E(A | L)} \\ & - \frac{(1 - A)\{Y - E(Y | A = 0, L)\}}{1 - E(A | L)} \\ & + E(Y | A = 1, L) - E(Y | A = 0, L) - \psi^*, \end{aligned}$$

is strictly larger than the variance of the locally efficient G-estimator (unless A and L are independent, as may be the case when A refers to a randomized treatment, in which case they are equally efficient). In particular, the asymptotic variance of the locally efficient IPW estimator equals

$$(22) \quad \sigma^2 E\left\{\frac{1}{\text{Var}(A | L)}\right\},$$

when the treatment effect is homogeneous. The difference between (21) and (22) can be sizeable when the propensity score is close to zero or 1 for some values of L for then $\text{Var}(A | L)$ is close to zero and thus $1/\text{Var}(A | L)$ can take on large values. In our opinion, this difference is not usually offset by the weaker restrictions imposed by the MSM. Indeed, the marginal treatment effect would seldom be of scientific interest when certain subjects are almost precluded from

receiving treatment or no treatment. Moreover, the G-estimator retains a useful interpretation even when the assumption of constant treatment effects fails in the sense that $E(Y^1 - Y^0 | A = 1, L) = \psi(L)$ for some function $\psi(L)$. Indeed, in that case the locally efficient G-estimator converges to

$$(23) \quad \frac{E\{\text{Var}(A | L)\psi(L)\}}{E\{\text{Var}(A | L)\}},$$

which continues to be useful as a weighted average of treatment effects $\psi(L)$, with most weight given to strata with most information about the treatment effect.

This difference in asymptotic variance between both estimators becomes even more pronounced in the likely event that the model for $E(Y | A = 0, L) = E(Y^0 | L) = E(Y - \psi^* A | L)$ is misspecified. Let $\Delta(L) = E(Y | A = 0, L) - E^*(Y | A = 0, L)$ denote the degree of misspecification at covariate value L , with $E(Y | A = 0, L)$ the true expectation and $E^*(Y | A = 0, L)$ the expectation used for evaluating the locally efficient G-estimator. Furthermore, assume that in truth the treatment effect is homogeneous. Then the asymptotic variance of the G-estimator becomes

$$\frac{\sigma^2}{E\{\text{Var}(A | L)\}} + \frac{E\{\text{Var}(A | L)\Delta(L)^2\}}{E\{\text{Var}(A | L)\}^2},$$

and the asymptotic variance of the previously considered IPW estimator becomes

$$\sigma^2 E\left\{\frac{1}{\text{Var}(A | L)}\right\} + E\left[\frac{\{\Delta(L) + \psi^* E(1 - A | L)\}^2}{\text{Var}(A | L)}\right].$$

Consider now that model misspecification is more likely in regions of little overlap. Then because $\text{Var}(A | L) \approx 0$ in these regions, model misspecification in these regions will only have a minor impact on the variance of the G-estimator, but a particularly strong impact on the variance of the locally efficient IPW estimator. Similar findings have been noted concerning the asymptotic bias of these estimators (Vansteelandt, Bekaert and Claeskens, 2012).

While this contrast between G-estimation and IPW-estimation under misspecification of the outcome model could turn out to be somewhat less dramatic when the propensity score is not considered as fixed and known, we believe that the above findings more likely understate the factual differences if one considers that mainstream applications are based on sequential treatments (and thus even more variable inverse probability

weights) and on simple, inefficient inverse probability weighting methods. The latter can be viewed as inducing extreme misspecification in the outcome model as they amount to setting $E(Y | A = 1, L) = E(Y | A = 0, L) = 0$. We thus believe that more routine application of G-estimation is warranted.

5. STRUCTURAL NESTED MODELS FOR TIME-VARYING TREATMENTS

Before introducing SNMs for time-varying or sequential treatments, we consider the structure of observed data in observational studies with repeated treatments and covariates, as well as definitions of causal effects in such setting. Suppose that measurements are collected at fixed time points t_0, t_1, \dots, t_{K+1} . Let A_k denote the treatment provided at time $t_k, k = 0, \dots, K$, and L_k denote other covariates measured at that time; Y_k , the outcome measured at time $t_k, k = 1, \dots, K + 1$, is part of L_k . We presume the variables are ordered L_0, A_0, L_1, A_1 , etc.; thus, covariates and outcome at t_k precede treatment at t_k .

Let $Y_m^{\bar{a}_{m-1}}$ denote the outcome that would be seen at time t_m in a given individual were (s)he to receive treatment history \bar{a}_{m-1} through time t_{m-1} . The variables $Y_m^{\bar{a}_{m-1}}$ are potential outcomes, which are again linked to the observed data via the consistency assumption that $Y_m = Y_m^{\bar{a}_{m-1}}$ if $\bar{A}_{m-1} = \bar{a}_{m-1}$. We presume that treatment at or after t_m cannot affect outcome at times up to t_m ; thus, $Y_m^{\bar{a}_{m-1}, \bar{a}_m} = Y_m^{\bar{a}_{m-1}, \bar{a}_m^\dagger}$ for $\bar{a}_m \neq \bar{a}_m^\dagger$. Causal effects can now be defined as comparisons of potential outcomes $Y^{\bar{a}_K}$ for the same group of subjects for different treatment histories $\bar{a}_K, \bar{a}_K^\dagger, \bar{a}_K \neq \bar{a}_K^\dagger$ (Robins, 1986). If the outcome is measured only at the end of a fixed follow-up period, or only at a subset of the follow-up times, we can let $Y_m = (\cdot)$, where “ \cdot ” denotes missing or undefined values for the times where the outcome is not measured. Most of the subsequent presentation then applies to those settings.

5.1 Structural Nested Mean Models

Structural nested mean models (SNMMs) (Robins, 1994; Robins, Rotnitzky and Scharfstein, 2000) simulate the sequential removal of an amount (“blip”) of treatment at t_m on subsequent average outcomes, after having removed the effects of all subsequent treatments. Given a history \bar{a}_m , define the counterfactual history $(\bar{a}_m, 0)$ as the history \bar{a}^\dagger that agrees with \bar{a}_m through time t_m and is 0 thereafter. SNMMs then

model the effect of a blip of treatment at t_m on the subsequent outcome means when holding all future treatments fixed at their reference level 0; thus, they parameterize contrasts of $\underline{Y}_{m+1}^{\bar{a}_m,0}$ and $\underline{Y}_{m+1}^{\bar{a}_{m-1},0}$ conditionally on treatment and covariate histories through t_m as

$$\begin{aligned} & g\{E(\underline{Y}_{m+1}^{\bar{a}_m,0} \mid \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m)\} \\ & - g\{E(\underline{Y}_{m+1}^{\bar{a}_{m-1},0} \mid \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m)\} \\ & = \gamma_m^*(\bar{l}_m, \bar{a}_m; \psi^*), \end{aligned}$$

for each $m = 0, \dots, K$ and (\bar{l}_m, \bar{a}_m) , where $\gamma_m^*(\bar{l}_m, \bar{a}_m; \psi)$ is a known $(K + 1 - m)$ -dimensional function, smooth in ψ , and for each \bar{l}_m, \bar{a}_{m-1} and ψ it is by definition required that $\gamma_m^*(\bar{l}_m, \bar{a}_{m-1}, 0; \psi) = 0$. Alternatively, one may focus on the effect of treatment on the end-of-study outcome $Y \equiv Y_{K+1}$ only, in which case one obtains a SNMM of the form

$$\begin{aligned} & g\{E(Y^{\bar{a}_m,0} \mid \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m)\} \\ & - g\{E(Y^{\bar{a}_{m-1},0} \mid \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m)\} \\ & = \gamma_m^*(\bar{l}_m, \bar{a}_m; \psi^*), \end{aligned}$$

for each $m = 0, \dots, K$ and (\bar{l}_m, \bar{a}_m) , where $\gamma_m^*(\bar{l}_m, \bar{a}_m; \psi)$ is now 1-dimensional. The above contrasts generalize the notion of the effect of treatment on the treated to the setting of a sequence of treatments. The name “nested” refers to the nesting across time, of the subgroups defined by \bar{L}_m and \bar{A}_m within which the effects are evaluated.

Typically, the parameterization is chosen to be such that $\gamma_m^*(\bar{l}_m, \bar{a}_m; 0) = 0$ for all \bar{l}_m, \bar{a}_m so that $\psi = 0$ encodes the null hypothesis of no treatment effect. For instance, with 2 time points ($K = 1$) a linear SNMM may be given by

$$\begin{aligned} & E(Y_2^{(a_0, a_1)} - Y_2^{(a_0, 0)} \mid \bar{L}_1 = \bar{l}_1, \bar{A}_1 = \bar{a}_1) \\ & = (\psi_0^* + \psi_1^* l_1 + \psi_2^* a_0) a_1, \\ & E(Y_2^{(a_0, 0)} - Y_2^0 \mid L_0 = l_0, A_0 = a_0) \\ & = (\psi_3^* + \psi_4^* l_0) a_0, \\ & E(Y_1^{(a_0, 0)} - Y_1^0 \mid L_0 = l_0, A_0 = a_0) \\ & = (\psi_0^* + \psi_1^* l_0) a_0. \end{aligned}$$

Here, the first equation models the effect of A_1 on Y_2 , the second models the effect of A_0 on Y_2 and the third models the effect of A_0 on Y_1 , all within levels of variables defined prior to the considered exposure. Thus, ψ_0^*, ψ_1^* and ψ_2^* encode short-term treatment effects, which are here assumed to be constant at all time

points, and ψ_3^* and ψ_4^* encode long-term treatment effects. These effects are visualised in Figures 2 and 3 below. When interest merely lies in the effect on the end-of-study outcome, then the above model for Y_1 can be ignored.

Under the SNMM, as in Section 2.1, it is possible to define a transformation $U_m^*(\psi^*)$ of \underline{Y}_{m+1} , whose mean value equals the mean that would be observed if treatment were suspended from time t_m onward, in the sense that

$$\begin{aligned} (24) \quad & E\{U_m^*(\psi^*) \mid \bar{L}_m, \bar{A}_{m-1} = \bar{a}_{m-1}, A_m\} \\ & = E(\underline{Y}_{m+1}^{\bar{a}_{m-1},0} \mid \bar{L}_m, \bar{A}_{m-1} = \bar{a}_{m-1}, A_m), \end{aligned}$$

for $m = 0, \dots, K$. Here, $U_m^*(\psi)$ is a vector with components

$$Y_k - \sum_{l=m}^{k-1} \gamma_{l,k}^*(\bar{L}_l, \bar{A}_l; \psi),$$

for $k = m + 1, \dots, K + 1$ (or for $k = K + 1$ only if interest merely lies in the effect on the end-of-study outcome) if $g(\cdot)$ is the identity link, and

$$Y_k \exp \left\{ - \sum_{l=m}^{k-1} \gamma_{l,k}^*(\bar{L}_l, \bar{A}_l; \psi) \right\},$$

if $g(\cdot)$ is the log link. These equations formalize the notion of peeling off or blipping down the treatment effects over the treatment period from t_m to t_{k-1} . For instance, in the previous example for 2 time points,

$$\begin{aligned} U_1^*(\psi^*) &= Y_2 - (\psi_0^* + \psi_1^* L_1 + \psi_2^* A_0) A_1, \\ U_0^*(\psi^*) &= (Y_1 - (\psi_0^* + \psi_1^* L_0) A_0, Y_2 \\ &\quad - (\psi_0^* + \psi_1^* L_1 + \psi_2^* A_0) A_1 \\ &\quad - (\psi_3^* + \psi_4^* L_0) A_0)'. \end{aligned}$$

For link functions other than the identity and log link, such a transformation can still be defined, but depends on the observed data distribution in a complicated and contrived way. For instance, when $g(\cdot)$ is the logit link and there are 2 time points ($K = 1$), then under the SNMM we have that

$$\begin{aligned} & E(Y_2^0 \mid L_0 = l_0, A_0 = a_0) \\ & = g^{-1} [g\{E(g^{-1}[g\{E(Y_2 \mid \bar{L}_1, A_1, A_0 = a_0)\} \\ &\quad - \gamma_1^*(\bar{L}_1, A_1, A_0 = a_0; \psi^*)] \\ &\quad \mid L_0 = l_0, A_0 = a_0)\} \\ &\quad - \gamma_0^*(l_0, a_0; \psi^*)]. \end{aligned}$$

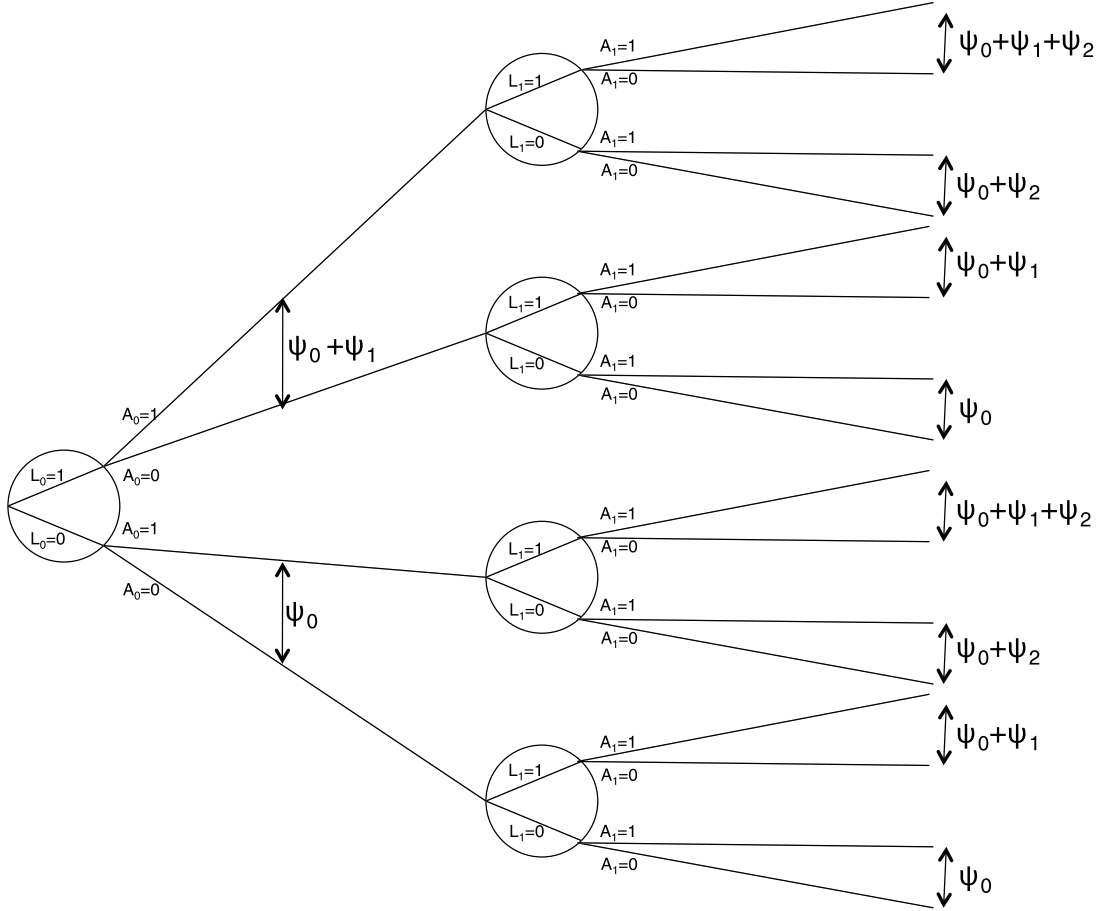


FIG. 2. Visualisation of the effects $E(Y_1^{(a_0,0)} - Y_1^0 | L_0 = l_0, A_0 = a_0)$ and $E(Y_2^{(a_0,a_1)} - Y_2^{(a_0,0)} | \bar{L}_1 = \bar{l}_1, \bar{A}_1 = \bar{a}_1)$. Lines within the circles depict covariate strata; lines outside the circles depict exposure strata.

The calculation of $U_0(\psi)$ thus not only demands knowledge of $E(Y_2 | \bar{L}_1, \bar{A}_1)$, but also of the distribution of (L_1, A_1) , given (L_0, A_0) .

The effect of a sequential treatment on the failure time distribution can be parameterized through a collection of SNMMs with log link, one for each time point (Robins and Hernan, 2009; Picciotto et al., 2012). In continuous time (Martinussen et al., 2011), such structural nested cumulative failure time models are defined by restrictions of the form:

$$\frac{P(T^{\bar{a}_m,0} > t | \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m, T \geq t_m)}{P(T^{\bar{a}_{m-1},0} > t | \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m, T \geq t_m)} = \exp\{\gamma_m^*(t, \bar{l}_m, \bar{a}_m; \psi^*)\},$$

for all t and $m = 0, \dots, K$, where $\gamma_m^*(t, \bar{l}_m, \bar{a}_m; \psi)$ is a known function, smooth in ψ and monotonic in t , and $\gamma_m^*(t, \bar{l}_m, \bar{a}_{m-1}, 0; \psi) = 0$ for all $t, \bar{l}_m, \bar{a}_{m-1}$ and ψ .

5.2 Structural Nested Distribution Models

Structural nested distribution models (SNDMs) are closely related to SNMMs, but parameterize a map between percentiles of the distribution of $Y_k^{\bar{a}_m,0}$ and percentiles of the distribution of $Y_k^{\bar{a}_{m-1},0}$. They are most easily understood by first considering the class of more restrictive rank-preserving SNDMs. In particular, for each exposure $A_m, m = 0, \dots, K$, let us first consider a rank-preserving SNDM to parameterize its effect on the end-of-study outcome Y :

$$Y^{\bar{a}_{m-1},0} = \gamma_m(Y^{\bar{a}_m,0}, \bar{l}_m, \bar{a}_m; \psi^*),$$

for subjects with $\bar{A}_m = \bar{a}_m$ and $\bar{L}_m = \bar{l}_m$, $m = 0, \dots, K$. Here, $\gamma_m(y, \bar{l}_m, \bar{a}_m; \psi)$ is a known function, smooth in ψ and a smooth, monotonic function of y , which contrasts the counterfactuals $Y^{\bar{a}_{m-1},0}$ and $Y^{\bar{a}_m,0}$, and must satisfy $\gamma_m(y, \bar{l}_m, \bar{a}_{m-1}, 0; \psi) = y$ for all y and ψ . For instance, with 2 time points ($K = 1$) a rank preserving SNDM may be given by the following set

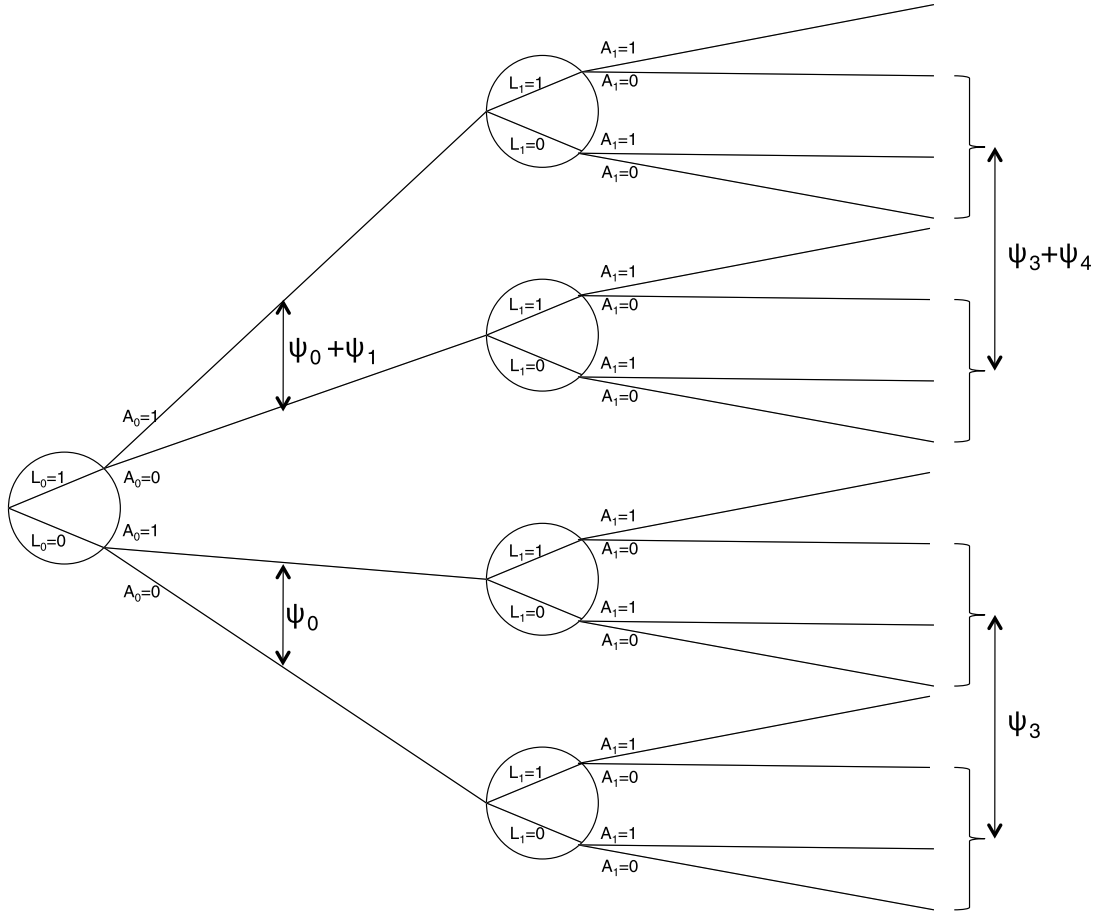


FIG. 3. Visualisation of the effects $E(Y_2^{(a_0,0)} - Y_2^0 | L_0 = l_0, A_0 = a_0)$. Lines within the circles depict covariate strata; lines outside the circles depict exposure strata.

of restrictions:

$$\begin{aligned} Y^{A_0,0} &= \gamma_1(Y, \bar{L}_1, \bar{A}_1) = Y - (\psi_1^* + \psi_2^* L_1) A_1, \\ Y^0 &= \gamma_0(Y^{A_0,0}, L_0, A_0) \\ &= Y^{A_0,0} - (\psi_1^* + \psi_2^* L_0) A_0 \\ &= Y - \psi_1^* (A_0 + A_1) - \psi_2^* (L_1 A_1 + L_0 A_0). \end{aligned}$$

A SNDM relaxes these restrictions by demanding that they merely hold in distribution, conditional on the observed history (i.e., $\bar{L}_m = \bar{l}_m$ and $\bar{A}_m = \bar{a}_m$).

To describe the effect on a repeated counterfactual future $Y_{m+k}^{\bar{a}_m,0}$, we can borrow ideas from Section 2.3.2. In particular, upon substituting A by A_m , L by \bar{L}_m , \bar{A}_{m-1} and Y_k by $Y_{m+k}^{\bar{a}_m,0}$ in the rank-preserving model (9), we obtain the identity:

$$(25) \quad Y_{m+k}^{\bar{a}_{m-1},0} = \gamma_{m,m+k}(Y_{m+1:m+k}^{\bar{a}_m,0}, \bar{l}_m, \bar{a}_m; \psi^*),$$

for subjects with $\bar{A}_m = \bar{a}_m$ and $\bar{L}_m = \bar{l}_m$, $m = 0, \dots, K$ and $k = 1, \dots, K + 1 - m$. Here, $\gamma_{m,m+k}(y_{m:m+k}, \bar{l}_m,$

$\bar{a}_m; \psi)$ is a known function, smooth in ψ and a smooth, monotonic function of y_{m+k} , which contrasts the counterfactuals $Y_{m+k}^{\bar{a}_{m-1},0}$ and $Y_{m+k}^{\bar{a}_m,0}$, and must satisfy $\gamma_{m,m+k}(y_{m:m+k}, \bar{l}_m, \bar{a}_{m-1}, 0; \psi) = y_{m+k}$ for all $y_{m+k}, \bar{l}_m, \bar{a}_{m-1}$ and ψ . For instance, with 2 time points ($K = 1$) a rank preserving SNDM may be given by the following set of restrictions:

$$\begin{aligned} (26) \quad Y_1^0 &= \gamma_{0,1}(Y_1, L_0, A_0; \psi^*) \\ &= Y_1 - (\psi_1^* + \psi_2^* L_0) A_0, \\ Y_2^{A_0,0} &= \gamma_{1,2}(Y_2, \bar{L}_1, \bar{A}_1; \psi^*) \\ &= Y_2 - (\psi_1^* + \psi_2^* L_1) A_1, \\ Y_2^0 &= \gamma_{0,2}(Y_1, Y_2^{A_0,0}, L_0, A_0; \psi^*) \\ &= Y_2^{(A_0,0)} - (\psi_3^* + \psi_4^* Y_1) A_0 \\ &= Y_2 - (\psi_1^* + \psi_2^* L_1) A_1 \\ &\quad - (\psi_3^* + \psi_4^* Y_1) A_0. \end{aligned}$$

Here, the first two equations express short-term exposure effects, that is, the effect of A_0 on Y_1 and of A_1 on Y_2 . The third equation expresses the effect of A_0 on Y_2 (more precisely, its effect on $Y_2^{A_0,0}$). As in Section 2.3.2, this equation must take into account that the effect may be different depending on the outcome level at time t_1 ; this allows for A_0 to also affect the dependence between Y_1 and Y_2 , but evidently complicates interpretation. More generally, rank-preserving SNDMs allow for the effect of a_m on Y_{m+k} , as encoded by a contrast of $Y_{m+k}^{\bar{a}_m,0}$ and $Y_{m+k}^{\bar{a}_{m-1},0}$, to depend on the history of treatments and covariates up to time t_m , but additionally on the potential outcome history under the treatment regime $(\bar{a}_m, 0)$, up to time t_{m+k-1} .

A SNDM relaxes the restrictions of a rank preserving SNDM by demanding that the equality (25) merely holds in distribution, conditional on $\bar{L}_m = \bar{l}_m$ and $\bar{A}_m = \bar{a}_m$. Assuming that for given (\bar{L}_m, \bar{A}_m) , \underline{Y}_{m+1} has a continuous multivariate distribution with probability 1, a SNDM can thus be defined by

$$(28) \quad \begin{aligned} F_{\underline{Y}_{m+1}^{\bar{a}_{m-1},0} | \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m} \{ \gamma_m(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_m; \psi^*) \} \\ = F_{\underline{Y}_{m+1}^{\bar{a}_m,0} | \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m}(\underline{y}_{m+1}), \end{aligned}$$

for all \bar{l}_m, \bar{a}_m , where $\gamma_m(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_m; \psi^*)$ is a vector with components $\gamma_{m,k}(\underline{y}_{m+1:m+k}, \bar{l}_m, \bar{a}_m; \psi^*)$ for $k = 1, \dots, K + 1 - m$, where the components $\gamma_{m,k}$ are defined in recursive fashion similar to in Section 2.3.2.

Under the SNDM, a variable $U_m(\psi^*) = (U_{m,m+1}(\psi^*), \dots, U_{m,K+1}(\psi^*))'$ can be constructed which predicts how the outcomes past time t_m would look like if treatment were suspended from time t_m onward, in the sense that

$$(29) \quad \begin{aligned} P\{U_m(\psi^*) > \underline{y}_{m+1} | \bar{L}_m, \bar{A}_m = \bar{a}_m\} \\ = P(\underline{Y}_{m+1}^{\bar{a}_{m-1},0} > \underline{y}_{m+1} | \bar{L}_m, \bar{A}_m = \bar{a}_m). \end{aligned}$$

This variable can be recursively obtained for $m = K, \dots, 0$ from

$$(30) \quad \begin{aligned} U_{m,m+k}(\psi) \\ \equiv \gamma_{m,m+k}\{(Y_{m+1}, U_{m+1:m+k}(\psi)), \bar{L}_m, \bar{A}_m; \psi\}, \end{aligned}$$

for $k = 1, \dots, K + 1 - m$, where we define $U_{m+1,m+k}(\psi)$ to be empty for $k = 1$. For instance, in the SNDM that assumes the identities in (26) hold in distribution (conditional on the observed history), we

have that

$$\begin{aligned} U_1(\psi) &= U_{1,2}(\psi) = \gamma_{1,2}(Y_2, \bar{L}_1, \bar{A}_1; \psi) \\ &= Y_2 - (\psi_1 + \psi_2 L_1) A_1, \\ U_0(\psi) &= (U_{0,1}(\psi), U_{0,2}(\psi)) \\ &= (\gamma_{0,1}(Y_1, L_0, A_0; \psi), \\ &\quad \gamma_{0,2}(Y_1, U_{1,2}(\psi), L_0, A_0; \psi)) \\ &= (Y_1 - (\psi_1 + \psi_2 L_0) A_0, \\ &\quad Y_2 - (\psi_1 + \psi_2 L_1) A_1 - (\psi_3 + \psi_4 Y_1) A_0). \end{aligned}$$

The identity (30) will be useful in estimation and for predicting the effect of specific interventions on the outcome distribution.

Structural nested failure time models (SNFTMs) are a variant of SNDMs which have seen most applications to date. These link percentiles from the conditional distributions of $T^{\bar{a}_{m-1},0}$ and $T^{\bar{a}_m,0}$, conditional on $\bar{L}_m, \bar{A}_m = \bar{a}_m$, and for subjects who are still in the risk set (say, alive) at time t_m :

$$\begin{aligned} S_{T^{\bar{a}_{m-1},0} | \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m, T \geq t_m} \{ \gamma_m(t, \bar{l}_m, \bar{a}_m; \psi^*) \} \\ = S_{T^{\bar{a}_m,0} | \bar{L}_m = \bar{l}_m, \bar{A}_m = \bar{a}_m, T \geq t_m}(t), \end{aligned}$$

for $t > t_m$, where $S(\cdot)$ denotes a survival function. Here, $\gamma_m(t, \bar{l}_m, \bar{a}_m; \psi^*)$ is a known function, smooth in ψ and monotonic in t , and $\gamma_m(t, \bar{l}_m, \bar{a}_{m-1}, 0; \psi) = t$ for all $t, \bar{l}_m, \bar{a}_{m-1}$ and ψ . For instance, the choice $\gamma_m(t, \bar{l}_m, \bar{a}_m; \psi) = t_m + (t - t_m) \exp(a_m \psi)$ for $t > t_m$ expresses that the effect of suspending treatment a_m at time t_m is to change the residual lifetime $t - t_m$ with a factor $\exp(a_m \psi)$. For this choice of model, one can predict among individuals who survive to (or through, or until) time t_m what their lifetime would be had treatment been suspended from time t_m onward, as

$$\begin{aligned} U_m(\psi) &= t_m + \sum_{k: t_m \leq t_k \leq T} (t_k - t_{k-1}) \exp(A_k \psi) \\ &\quad + (T - t_{T-}) \exp(A_{t_{T-}} \psi), \end{aligned}$$

where t_{T-} denotes the largest time point in $\{t_0, \dots, t_K\}$ less than T and $U_m(\psi)$ is a random variable for which (for $t > t_m$)

$$\begin{aligned} P\{U_m(\psi^*) > t | \bar{L}_m, \bar{A}_m = \bar{a}_m, T \geq t_m\} \\ = P(T^{\bar{a}_{m-1},0} > t | \bar{L}_m, \bar{A}_m = \bar{a}_m, T \geq t_m). \end{aligned}$$

5.3 Retrospective Blip Models

Retrospective blip models have been extended to model the effect of a sequential treatment on a scalar end-of-study outcome $Y \equiv Y_{K+1}$ conditional on the

treatment and covariate history up to end-of-study. Mean models take the form

$$(31) \quad \begin{aligned} & g\{E(Y^{\bar{a}_m,0} | \bar{L}_K = \bar{l}_K, \bar{A}_K = \bar{a}_K)\} \\ & - g\{E(Y^{\bar{a}_{m-1},0} | \bar{L}_K = \bar{l}_K, \bar{A}_K = \bar{a}_K)\} \\ & = \gamma_m^*(\bar{l}_K, \bar{a}_K; \psi^*), \end{aligned}$$

where $\gamma_m^*(\bar{l}_K, \bar{a}_K; \psi)$ is a known function, smooth in ψ and equaling zero for all ψ, \bar{l}_K and \bar{a}_K with $a_m = 0$. Distribution models take the form:

$$\begin{aligned} & F_{Y^{\bar{a}_{m-1},0} | \bar{L}_K = \bar{l}_K, \bar{A}_K = \bar{a}_K} \{\gamma_m(y, \bar{l}_K, \bar{a}_K; \psi^*)\} \\ & = F_{Y^{\bar{a}_m,0} | \bar{L}_K = \bar{l}_K, \bar{A}_K = \bar{a}_K}(y), \end{aligned}$$

where $\gamma_m(y, \bar{l}_K, \bar{a}_K; \psi)$ is a known function, smooth in ψ and equaling y for all ψ, y, \bar{l}_K and \bar{a}_K with $a_m = 0$; a rank-preserving version of this was proposed by Joffe, Small and Hsu (2007). For nonparametric identifiability, restrictions are needed on the functions $\gamma_m^*(\bar{l}_K, \bar{a}_K; \psi^*)$ and $\gamma_m(y, \bar{l}_K, \bar{a}_K; \psi)$, for example, that they do not involve the future \bar{a}_{m+1} and \bar{l}_{m+1} (Vansteelandt, 2010).

Retrospective blip models can be useful for modeling a dichotomous outcome (Vansteelandt, 2010). Under these models, identity (24) is satisfied with $U_m^*(\psi)$ being a vector with components

$$g^{-1} \left[g\{E(Y | \bar{L}_K, \bar{A}_K)\} - \sum_{l=m}^K \gamma_l^*(\bar{L}_K, \bar{A}_K; \psi) \right].$$

Evaluation of $U_m^*(\psi)$ (which is needed to make estimation of ψ^* manageable) then merely requires a model for $E(Y | \bar{L}_K, \bar{A}_K)$, but not for the distribution of treatment and covariates at each time. The parameters indexing these models are nonetheless more limited than the parameters indexing SNMMs in that they cannot be used by themselves for making treatment decisions prior to the end-of-study time, unless one integrates over the distribution of covariates subsequent to m (see, e.g., Vansteelandt, 2010).

6. IDENTIFICATION AND ESTIMATION IN STRUCTURAL NESTED MODELS FOR SEQUENTIAL TREATMENTS

This section sketches identifying assumptions and inferential methods for sequential treatments. Under instrumental variables assumptions sketched in Section 6.3 and under the future ignorability assumptions sketched in Section 6.2, inferential methods have been developed for SNMs, but these assumptions do not suffice for the identification of marginal treatment effects,

and hence parameters indexing MSMs. The broader array of useful identifying assumptions thus constitutes an important advantage of SNMs.

6.1 Sequential Ignorability

The assumption of ignorable treatment assignment can be generalised to sequential treatments as follows:

$$(32) \quad A_m \perp\!\!\!\perp Y_{m+1}^{\bar{a}_{m-1},0} | \bar{L}_m, \bar{A}_{m-1} = \bar{a}_{m-1},$$

for $m = 0, \dots, K$. This assumption has been called variously “no unmeasured confounders assumption,” “sequential ignorability,” “sequential randomization” or “exchangeability.” It expresses that at each time t_m , the observed history of covariates \bar{L}_m and exposures \bar{A}_{m-1} includes all risk factors of A_m that are also associated with future outcomes.

This assumption together with identity (24) imply that

$$E\{U_m(\psi^*) | \bar{L}_m, \bar{A}_m\} = E\{U_m(\psi^*) | \bar{L}_m, \bar{A}_{m-1}\}$$

for all m under a SNMM. The parameter ψ^* indexing a SNMM can therefore be estimated by solving

$$(33) \quad \begin{aligned} 0 = & \sum_{i=1}^n \sum_{m=0}^K [d_m(\bar{L}_{im}, \bar{A}_{im}) \\ & - E\{d_m(\bar{L}_{im}, \bar{A}_{im}) | \bar{L}_{im}, \bar{A}_{i,m-1}\}] \\ & \times [U_{im}(\psi) - E\{U_{im}(\psi) | \bar{L}_{im}, \bar{A}_{i,m-1}\}], \end{aligned}$$

where $d_m(\bar{L}_{im}, \bar{A}_{im})$, $m = 0, \dots, K$ is an arbitrary $p \times (K + 1 - m)$ -dimensional function, with p the dimension of ψ . This estimating equation essentially sets the sum across time points m of the conditional covariances between $U_{im}(\psi)$ and the given function $d_m(\bar{L}_{im}, \bar{A}_{im})$, given $\bar{L}_{im}, \bar{A}_{i,m-1}$, to zero. When the previous outcome is included in the confounder history (i.e., \bar{L}_{im} includes Y_{im}) and there is homoscedasticity [i.e., when the conditional variance of $U_{im}(\psi^*)$ given $\bar{L}_{im}, \bar{A}_{im}$ is constant for $m = 0, \dots, K$], then local semiparametric efficiency under the SNMM is attained upon choosing

$$d_m(\bar{L}_{im}, \bar{A}_{im}) = E \left\{ \frac{\partial U_m(\psi^*)}{\partial \psi} | \bar{L}_{im}, \bar{A}_{im} \right\}.$$

Sequential ignorability (32) together with identity (29) moreover implies that

$$U_m(\psi^*) \perp\!\!\!\perp A_m | \bar{L}_m, \bar{A}_{m-1}$$

for all m under the SNMM. This conditional independence restriction suggests that the parameter indexing

a SNDM can be solved from

$$\begin{aligned}
 0 &= \sum_{i=1}^n \sum_{m=0}^K d_m \{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\} \\
 &\quad - E[d_m \{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\} | \bar{L}_{im}, \bar{A}_{im}] \\
 (34) \quad &- E[d_m \{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\} \\
 &\quad - E[d_m \{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\} | \bar{L}_{im}, \bar{A}_{im}] | \\
 &\quad U_{im}(\psi), \bar{L}_{im}, \bar{A}_{i,m-1}),
 \end{aligned}$$

where the index functions $d_m \{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\}$ must be of the dimension of ψ . When the previous outcome is included in the confounder history (i.e., \bar{L}_{im} includes Y_{im}), then local semiparametric efficiency is obtained upon choosing

$$\begin{aligned}
 d_m \{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\} \\
 = E\{S_\psi(\psi) | U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\},
 \end{aligned}$$

where $S_\psi(\psi)$ is the score for ψ under the observed data likelihood

$$\begin{aligned}
 f(\bar{Y}_{K+1}, \bar{L}_K, \bar{A}_K) \\
 = f\{U_0(\psi^*)\} \\
 \cdot \prod_{m=0}^K \left[f\{L_m | \bar{L}_{m-1}, \bar{A}_{m-1}, U_m(\psi^*)\} \right. \\
 \times f\{A_m | \bar{L}_m, \bar{A}_{m-1}, U_m(\psi^*)\} \\
 \left. \cdot \left| \frac{\partial U_m(\psi^*)}{\partial U_{m+1}(\psi^*)} \right| \right],
 \end{aligned}$$

with all components substituted by suitable parametric models (Robins, 1997); here, the term $f\{A_m | \bar{L}_m, \bar{A}_{m-1}, U_m(\psi^*)\} = f(A_m | \bar{L}_m, \bar{A}_{m-1})$ under sequential ignorability, and thus can be ignored. This likelihood formulation is of interest in itself because it enables specifying the joint distribution of the variables in a way that is consistent with the sharp null hypothesis of no effect under the assumption of sequential ignorability, even in the presence of confounding by variables affected by treatment, which turns out more difficult with standard parameterisations (Robins, 1997).

Solving estimating equations (33) and (34) requires a parametric model \mathcal{A} for the conditional distribution of the exposure A_m for $m = 0, \dots, K$:

$$f(A_m | \bar{L}_{m-1}, \bar{A}_{m-1}) = f(A_m | \bar{L}_{m-1}, \bar{A}_{m-1}; \alpha^*),$$

where $f(A_m | \bar{L}_{m-1}, \bar{A}_{m-1}; \alpha)$ is a known density function, smooth in α , and α^* is an unknown finite-dimensional parameter which can be estimated via

standard maximum likelihood. In addition, it requires a parametric model \mathcal{B} for the conditional mean (or distribution) of $U_m^*(\psi^*)$ [or $U_m(\psi^*)$] for $m = 0, \dots, K$:

$$\begin{aligned}
 f\{U_m(\psi^*) | \bar{L}_m, \bar{A}_{m-1}\} \\
 = f\{U_m(\psi^*) | \bar{L}_m, \bar{A}_{m-1}; \gamma^*\},
 \end{aligned}$$

where $f\{U_m(\psi^*) | \bar{L}_m, \bar{A}_{m-1}; \gamma\}$ is a known density function, smooth in γ and γ^* is an unknown finite-dimensional parameter. As before, when the parameters α and γ are variation-independent, then so-called G-estimators that solve (33) and (34), obtained upon substituting α^* and γ^* by consistent estimators, are doubly robust (Robins and Rotnitzky, 2001): consistent when the SNM and either model \mathcal{A} or model \mathcal{B} is correctly specified, regardless of which. This double robustness property of the G-estimator is desirable for various reasons. First, it provides justification for using simple models for the multivariate distribution $f\{U_m(\psi^*) | \bar{L}_m, \bar{A}_{m-1}\}$ or even setting $E\{U_{im}(\psi) | \bar{L}_{im}, \bar{A}_{i,m-1}\} = 0$ in (33) for computational convenience. Second, while alternative proposals that rely on correct specification of model \mathcal{B} (see, e.g., Almirall, Ten Have and Murphy, 2010; Henderson, Ansell and Alshibani, 2010) tend to give more efficient estimators (under correct model specification), the concern for misspecification of model \mathcal{B} may be considerable in view of the aforementioned difficulty of postulating this model. This distribution can indeed be difficult to specify in view of its multivariate nature, the fact that $U_m(\psi^*)$ represents a transformation of the observed data and that it may moreover share the same outcome over multiple time points, so that the models for $U_m(\psi^*)$ corresponding to different time points may not be congenial at all times. This concern can be overcome by inferring the conditional expectations $E\{U_m(\psi^*) | \bar{L}_m, \bar{A}_{m-1}\}$ from models for the conditional distribution of L_{m+1} given \bar{L}_m, \bar{A}_{m-1} at each time m (Robins, Rotnitzky and Scharfstein, 2000; Almirall, Ten Have and Murphy, 2010). However, when the covariate \bar{L}_m is high-dimensional and/or strongly associated with treatment A_m , specifying such models can be a thorny and nontrivial task.

6.2 Departures from Sequential Ignorability and Sensitivity Analysis

Specified departures from (32) can also yield identification. For instance, one can allow dependence of treatment on a specified portion of the future potential outcomes, by relaxing (32) to (Joffe, Yang and Feldman, 2010; Zhang, Joffe and Small, 2011)

$$A_m \perp\!\!\!\perp Y_{m+\Delta}^{\bar{A}_{m-1},0} | \bar{L}_m, \bar{A}_{m-1} = \bar{a}_{m-1}, Y_{m+1:m+\Delta-1}^{\bar{A}_{m-1},0},$$

for some integer $\Delta \geq 1$; such assumptions have been termed future ignorability, since the independence at m is conditional on potential outcomes referring to times after m . This assumption can sometimes eliminate residual confounding bias, for instance, because the treatment process occurs in continuous time but confounding covariates are only measured intermittently, as is common in observational studies (Zhang, Joffe and Small, 2011), or when the future potential outcomes serve as proxies for other unmeasured confounding variables (Rosenbaum, 1984). However, it does not lead to nonparametric identification of the SNM parameters, so that inference becomes more dependent on correct specification of the causal model.

Alternatively, deviations from sequential ignorability can be parameterized as

$$\begin{aligned}
 & f(A_m = a_m \mid \bar{L}_m = \bar{l}_m, \\
 & \quad \bar{A}_{m-1} = \bar{a}_{m-1}, \underline{Y}_{m+1}^{\bar{a}_{m-1}, 0} = \underline{y}_{m+1}) \\
 & = t(A_m = a_m \mid \bar{L}_m = \bar{l}_m, \bar{A}_{m-1} = \bar{a}_{m-1}) \\
 & \quad \cdot \exp\{q_m(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_m)\} \\
 & \quad \cdot \left(\int t(A_m = a_m^\dagger \mid \bar{L}_m = \bar{l}_m, \bar{A}_{m-1} = \bar{a}_{m-1}) \right. \\
 & \quad \cdot \exp\{q_m(\underline{y}_{m+1}, \bar{l}_m, (\bar{a}_{m-1}, a_m^\dagger))\} da_m^\dagger \Big)^{-1},
 \end{aligned} \tag{35}$$

with $q_m(\cdot)$ known, satisfying $q_m(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_{m-1}, a_m = 0) = 0$ for all $(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_{m-1})$ and with $t(A_m \mid \bar{L}_m, \bar{A}_{m-1})$ an unknown conditional density. With $q_m(\cdot) = 0$ encoding the assumption of sequential ignorability, the function $q_m(\cdot)$ thus expresses the degree of departure from that assumption. As the data carry no genuine information about it, progress must be made by repeating the analysis with $q_m(\cdot)$ fixed at different values, which are then varied over some plausible range (Robins, Rotnitzky and Scharfstein, 2000); for example, by setting $q_m(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_m) = \gamma y_{m+1} a_m$, where γ is varied between -1 and 1 .

6.3 Instrumental Variables Assumptions

When the assumption of sequential ignorability fails, progress can sometimes be also made using an instrumental variable (IV). Such variable A_0 is assumed to satisfy

$$A_0 \perp\!\!\!\perp \underline{Y}^0 \mid L_0 \tag{36}$$

and

$$F_{\underline{Y}^0 \mid L_0=l_0, A_0=a_0}(\underline{y}) = F_{\underline{Y}^{a_0, 0} \mid L_0=l_0, A_0=a_0}(\underline{y}) \tag{37}$$

for all a_0, l_0 (Robins, 1989). Both these assumptions together imply that the instrument A_0 is not associated with the outcome, except through its association with subsequent treatments $A_m, m \geq 1$, which may affect outcome. These or similar assumptions have been used in adjusting for noncompliance in randomized trials (Robins and Tsiatis, 1991; Mark and Robins, 1993; Robins, 1994). With $A_m, m \geq 1$ denoting actual treatment and A_0 denoting randomized treatment, these assumptions are plausible when randomization does not affect the outcome other than by influencing the actual treatment.

Estimation under the IV assumptions can be based on estimating equations (33) and (34), but requires setting $d_m(\bar{L}_{im}, \bar{A}_{im}) = 0$ and $d_m\{U_{im}(\psi), \bar{A}_{im}, \bar{L}_{im}\} = 0$ for $m > 0$. Because of these restrictions, root- n estimation of ψ^* typically requires additional assumptions on $\gamma_m^*(\bar{l}_m, \bar{a}_m; \psi^*)$ and $\gamma_m(\underline{y}_{m+1}, \bar{l}_m, \bar{a}_m; \psi^*)$. In particular, it is commonly assumed that these functions are linear in a_m and do not involve a_0 ; moreover, time-varying covariates are commonly ignored, that is, L_m is set empty for $m > 0$. For instance, in linear SMMs for a single treatment A_1 (i.e., when $K = 1$) and dichotomous instrument, $\omega(L_0)$ in

$$E(Y_1 - Y_1^{a_0 0} \mid L_0 = l_0, \bar{A}_1 = \bar{a}_1) = \omega(l_0)a_1,$$

is just identified. Thus residual dependencies on a_0 or nonlinear dependencies on a_1 cannot be identified unless other untestable assumptions are imposed.

The resulting class of G-estimators contains the popular two-stage least squares estimator as a special case (Okui et al., 2012). However, the framework of G-estimation for SNMMs and SNDMs has the advantage that it extends immediately to outcomes that do not lend themselves to linear modeling, for example, censored failure-time outcomes (Robins and Tsiatis, 1991) and dichotomous outcomes (Vansteelandt and Goetghebeur, 2003; Robins and Rotnitzky, 2004), as well as to sequential treatments (Robins and Hernan, 2009). For instance, when $K = 1$ and L_1 is empty, the logistic SMM

$$\frac{\text{odds}(Y_2^{a_1} = 1 \mid L_0 = l_0, \bar{A}_1 = \bar{a}_1)}{\text{odds}(Y_2^0 = 1 \mid L_0 = l_0, \bar{A}_1 = \bar{a}_1)} = \exp(\psi^* a_1),$$

can be fitted by solving the SMM estimating equations with $U^*(\psi)$ given by $\text{expit}\{\text{logit } E(Y_2 \mid \bar{A}_1, L_0) - \psi A_1\}$ [cfr. (3)] and $E(Y_1 \mid \bar{A}_1, L_0)$ substituted by the fitted value under a parametric model (Vansteelandt and Goetghebeur, 2003; Vansteelandt et al., 2011). This additional model may sometimes not be congenial with the SMM and instrumental variables assumptions

in the sense that there may be no choice of parameter values indexing this model that satisfies these assumptions. This can be overcome by avoiding parameterization of the main effect of A_0 (conditional on L_0) in the model for $E(Y_1 | \bar{A}_1, L_0)$ and instead modeling the distribution of A_1 , given A_0 and L_0 (Robins and Rotnitzky, 2004), or by completely saturating the parameterization of the main effect of A_0 (conditional on L_0) (Vansteelandt et al., 2011). van der Laan, Hubbard and Jewell (2007) abandon logistic SMMs in favor of an interesting, but difficult to interpret relative risk parameterization. Alternatively, multiplicative SMMs can be used; under such models, case-only estimators have been constructed, which remain valid under case-control sampling (Bowden and Vansteelandt, 2011).

Variant assumptions have been proposed that allow use of time-varying instruments along with SNMMs and G-estimation. Robins and Hernan (2009) consider settings in which, at each time point, there is a variable whose association with the outcome of interest may be explained solely by its association with prior history and its effect on some treatment of interest. Joffe, Yang and Feldman (2010) consider settings in which the conditional independence of treatment and future potential outcomes in (32) holds for only an identifiable subset $\{i, m\}$ of the person-observations in the population rather than for all such observations. Treatment assignment in that subset may thus be considered an instrument for its effect and the effect of subsequent treatments.

IV analyses have several drawbacks relative to those based on sequential ignorability: (1) nonparametric identification is lost, and so inference is more dependent on correct specification of the causal model; (2) decreased power and precision; and (3) larger finite-sample bias.

6.4 Censoring

In SNFTMs, Type I censoring can be dealt with as previously explained by substituting $U_m(\psi)$ by an arbitrary function of $X_m(\psi) \equiv \min\{U_m(\psi), C_m(\psi)\}$ and $\Delta_m(\psi) \equiv I\{U_m(\psi) < C_m(\psi)\}$, where

$$C_m(\psi) \equiv \min\{U_m(C, \bar{a}_C, \bar{l}_C; \psi); \bar{a}_C, \bar{l}_C \in LA_m(C)\},$$

where $LA_m(C)$ is a given set of (\bar{a}_C, \bar{l}_C) histories which agree with the observed history of L through time t_m or C , whichever comes first, and A through time t_{m-1} or C , whichever comes first, and where $U_m(C, \bar{a}_C, \bar{l}_C; \psi)$ is defined like $U_m(\psi)$ in Section 5.2, but with C replacing T and \bar{a}_C and \bar{l}_C replacing \bar{A}_T and \bar{L}_T .

7. PREDICTING THE EFFECTS OF INTERVENTIONS

Identities (24) and (30) suggest using $U_m^*(\psi^*)$ and $U_m(\psi^*)$, respectively, as a prediction of $Y_{m+1}^{\bar{a}_{m-1}, 0}$ among individuals with observed history $\bar{A}_m = \bar{a}_m$. In particular, $E(\underline{Y}^0 | A_0, L_0) = E\{U_0^*(\psi^*) | A_0, L_0\}$ in SNMMs and $E(\underline{Y}^0 | A_0, L_0) = E\{U_0(\psi^*) | A_0, L_0\}$ in SNDMs, so that the expected outcome in the absence of treatment can be estimated as the sample average of $U_0^*(\hat{\psi})$ in SNMMs and of $U_0(\hat{\psi})$ in SNDMs. To estimate $E(\underline{Y}^{\bar{a}_K})$ for a different treatment regime \bar{a}_K , one could use a different structural nested model (SNM) with \bar{a}_K as the reference treatment regime. However, when—as often—the interest lies in comparing the expected counterfactual outcomes between different treatment regimes, then a concern is that these different SNMs may fail to imply a coherent model. Further complications arise when the goal is to evaluate the expected counterfactual outcome following a dynamic treatment regime whereby the treatment at each time t_m is assigned as a function of the treatment and covariate history up to that time; that is, for each m , $a_m = g(\bar{a}_{m-1}, \bar{l}_m)$.

These complications can be overcome by supplementing the SNM with so-called current treatment interaction functions (Robins, Rotnitzky and Scharfstein, 2000) about which the data carry no information, but which enable one to transport treatment effects in the treated to population-averaged treatment effects. For instance, let $K = 1$ and suppose that a SNMM has been fitted with $g(\cdot)$ the identity link. For simplicity, we focus here only on the effect of a nondynamic regime (a_0, a_1) at an end-of-study outcome $Y = Y_2$; results for dynamic treatment regimes are recovered upon making the substitutions $(g_0(\bar{l}_0), g_1(a_0, \bar{l}_1))$ for (a_0, a_1) . Two current treatment interaction functions can be defined, one for each sequential treatment:

$$\begin{aligned} r_1^*(\bar{L}_1, \bar{a}_1) &= E(Y^{a_0 a_1} - Y^{a_0 0} | A_0 = a_0, A_1 = a_1, \bar{L}_1) \\ &\quad - E(Y^{a_0 a_1} - Y^{a_0 0} | A_0 = a_0, A_1 \neq a_1, \bar{L}_1), \\ r_0^*(L_0, \bar{a}_1) &= E(Y^{a_0 a_1} - Y^0 | A_0 = a_0, L_0) \\ &\quad - E(Y^{a_0 a_1} - Y^0 | A_0 \neq a_0, L_0). \end{aligned}$$

These express how much the effects of subsequent treatment at m [i.e., a_1 and (a_0, a_1) at times 1 and 0, resp.] differ between groups that received that

level of treatment at m and those that did not. Under the SNMM, it is easily deduced from knowledge of $r_1^*(\bar{L}_1, \bar{a}_1)$ and $r_0^*(L_0, \bar{a}_1)$ that $E(Y^{a_0 a_1} - Y^0 \mid A_0 = a_0, L_0)$ equals

$$\begin{aligned} & E(Y^{a_0 a_1} - Y^{a_0 0} \mid A_0 = a_0, L_0) \\ & + E(Y^{a_0 0} - Y^0 \mid A_0 = a_0, L_0) \\ & = E\{\gamma_1^*(\bar{L}_1, \bar{a}_1; \psi^*) \\ & \quad - r_1^*(\bar{L}_1, \bar{a}_1)P(A_1 \neq a_1 \mid A_0 = a_0, \bar{L}_1) \mid \\ & \quad A_0 = a_0, L_0\} \\ & + \gamma_0^*(L_0, a_0; \psi^*). \end{aligned}$$

Because $E(Y^{a_0 a_1} - Y^0 \mid L_0)$ moreover equals

$$\begin{aligned} & E(Y^{a_0 a_1} - Y^0 \mid A_0 = a_0, L_0) \\ & - r_0^*(L_0, \bar{a}_1)P(A_0 \neq a_0 \mid L_0), \end{aligned}$$

we thus obtain that $E(Y^{a_0 a_1}) = E(Y^{a_0 a_1} - Y^0) + E(Y^0)$ equals

$$\begin{aligned} & E[E\{\gamma_1^*(\bar{L}_1, \bar{a}_1; \psi^*) \\ & \quad - r_1^*(\bar{L}_1, \bar{a}_1)P(A_1 \neq a_1 \mid A_0 = a_0, \bar{L}_1) \mid \\ & \quad A_0 = a_0, L_0\} \\ & + \gamma_0^*(L_0, a_0; \psi^*) \\ & - r_0^*(L_0, \bar{a}_1)P(A_0 \neq a_0 \mid L_0) + U_0^*(\psi^*)]. \end{aligned}$$

When there is no current treatment interaction [i.e., $r_1^*(\bar{l}_1, \bar{a}_1) = r_0^*(l_0, \bar{a}_1) = 0$ for all a_0, a_1, l_0, l_1], we thus have that

$$\begin{aligned} & E(Y^{a_0 a_1}) \\ & = E[E\{\gamma_1^*(\bar{L}_1, \bar{a}_1; \psi^*) \mid A_0 = a_0, L_0\} \\ & \quad + \gamma_0^*(L_0, a_0; \psi^*) + U_0^*(\psi^*)]. \end{aligned}$$

While the components $\gamma_1^*(\bar{L}_1, \bar{a}_1; \psi^*)$, $\gamma_0^*(L_0, a_0; \psi^*)$ and $U_0^*(\psi^*)$ can be estimated along the lines described in previous sections, a complication is that a model for the distribution of L_1 , conditional on A_0, L_0 , is needed to evaluate this; this can be cumbersome when L_1 is high-dimensional. This complication is avoided in simple structural models in which there is no effect modification by post-treatment variables [i.e., $\gamma_1^*(\bar{L}_1, \bar{a}_1)$ is not a function of L_1] and nondynamic regimes are considered.

The assumption of no current treatment interaction is satisfied under a mild strengthening of sequential ignorability such that

$$A_m \perp\!\!\!\perp \underline{Y}_{m+1}^{\bar{a}_K} \mid \bar{L}_m, \bar{A}_{m-1} = \bar{a}_{m-1},$$

for all m and all treatment histories \bar{a}_K . It is likewise sometimes satisfied under a mild strengthening of the instrumental variables assumption (36) such that for all treatment histories \bar{a}_K :

$$A_0 \perp\!\!\!\perp \underline{Y}^{\bar{a}_K} \mid L_0,$$

and a mild strengthening of the structural model such that, for instance, for binary A_1 (0/1):

$$\begin{aligned} & E(Y^{a_0 a_1} - Y^{a_0 a_1^\dagger} \mid A_1 = a_1, A_0 = a_0, L_0) \\ & = \gamma_1^*(a_1^\dagger, L_0; \psi^*)(a_1 - a_1^\dagger), \end{aligned}$$

for all a_1, a_1^\dagger . Following the instrumental variables assumptions, $Y^{a_0 0}$ and $Y^{a_0 1}$ should then be independent of A_0 , given L_0 , which respectively implies that

$$\begin{aligned} & E\{Y - \gamma_1^*(0, L_0; \psi^*)A_1 \mid A_0, L_0\} \\ & = E\{Y - \gamma_1^*(0, L_0; \psi^*)A_1 \mid L_0\}, \\ & E\{Y - \gamma_1^*(1, L_0; \psi^*)(1 - A_1) \mid A_0, L_0\} \\ & = E\{Y - \gamma_1^*(1, L_0; \psi^*)(1 - A_1) \mid L_0\}. \end{aligned}$$

It follows from this that $\gamma_1^*(0, L_0; \psi^*) = -\gamma_1^*(1, L_0; \psi^*)$, and thus again that the no current treatment interaction assumption is satisfied (Hernan and Robins, 2006).

8. DIRECT AND INDIRECT EFFECTS

SNMs parameterize the effects of treatment at each time with subsequent treatments set to some reference level. These effects can be viewed as controlled direct effects (Robins and Greenland, 1992), controlling all subsequent treatments at their reference levels. The formalism of SNMs is therefore more widely applicable for inferring the controlled direct effect of some target exposure A_0 on an outcome Y , other than through some mediator A_1 (e.g., the direct effect of the FTO gene on the risk of myocardial infarction other than via body mass). In particular, in the SNMM

$$\begin{aligned} & E(Y - Y^{a_0 0} \mid \bar{A}_1 = \bar{a}_1, \bar{L}_1) = \gamma_1^*(\bar{a}_1, \bar{L}_1; \psi^*), \\ & E(Y^{a_0 0} - Y^0 \mid A_0 = a_0, L_0) = \gamma_0^*(a_0, L_0; \psi^*), \end{aligned}$$

$\gamma_0^*(a_0, L_0; \psi^*)$ encodes the controlled direct effect of setting A_0 to zero, controlling A_1 at zero uniformly in the population. However, caution is warranted because $\gamma_0^*(a_0, L_0; \psi^*)$ may not encode the controlled direct effect of setting A_0 to zero, when controlling A_1 at some value $a_1 \neq 0$ (Robins and Wasserman, 1997). From knowledge that $\gamma_0^*(a_0, l_0; \psi^*)$ for all a_0, l_0 , one thus cannot deduce that A_0 has no direct effect on Y

(other than via A_1). Robins (1999) therefore proposed directly parameterizing the controlled direct effect as

$$(38) \quad \begin{aligned} E(Y^{a_0 a_1} - Y^{0 a_1} \mid A_0 = a_0, L_0) \\ = m(a_0, a_1, L_0; \psi^*), \end{aligned}$$

where $m(a_0, a_1, L_0; \psi)$ is a known function, smooth in ψ , which satisfies $m(0, a_1, l_0; \psi) = 0$. In contrast to SNMMs, (38) parameterizes only the effect of a_0 ; in (38), a_1 may, however, be a modifier of the effect of a_0 .

Since model (38) for fixed a_1 is a SMM for the counterfactual outcome Y^{a_1} , the techniques of Section 3 would be applicable to estimate ψ^* if Y^{a_1} were observed for each subject. Since Y^{a_1} is only observed for individuals with exposure level a_1 , Robins (1999) proposed treating subjects who receive a level of A_1 other than a_1 as censored and, assuming sequential ignorability, to inversely weight the data by the density $f(A_1 \mid A_0, \bar{L}_1)$ to control resulting selection bias. This amounts to solving ψ from an estimating equation of the form

$$(39) \quad \begin{aligned} 0 = & \frac{1}{n} \sum_{i=1}^n \frac{1}{f(A_{i1} \mid A_{i0}, \bar{L}_{i1})} [d(A_{i0}, L_{i0}) \\ & - E\{d(A_{i0}, L_{i0}) \mid L_{i0}\}] \\ & \times [Y_i - m(\bar{A}_{i1}, L_{i0}; \psi) \\ & - E\{Y_i - m(\bar{A}_{i1}, L_{i0}; \psi) \mid L_{i0}\}], \end{aligned}$$

where $d(A_{i0}, L_{i0})$ is an arbitrary index function. More efficient and doubly robust estimators have been reported elsewhere (Goetgeluk, Vansteelandt and Goetghebeur, 2008), as well as extensions to time-varying treatments (Robins, 1999).

Ignorability assumptions can be violated even in randomized trials (and Mendelian randomization studies), where assumption (36) is guaranteed by design, but the processes underlying the evolution of subsequent mediators may be poorly understood. Robins and Greenland (1994) avoid ignorability assumptions concerning the mediators by using initial randomization (or more generally, instrumental variables assumptions) to estimate controlled direct effects with SNFTMs. One can also use these approaches with SNMMs or SNDMs (e.g., Ten Have et al., 2007) and, in principle, in the presence of multiple mediators.

SMMs have also been developed for so-called natural direct effects (Robins and Greenland, 1992; Pearl, 2001). With $Y^{a_0 A_1^0}$ denoting the counterfactual outcome if A_0 were set to a_0 and A_1 to the counterfactual level A_1^0 that A_1 would take if A_0 were set to

zero, these are defined by contrasts between $Y^{a_0 A_1^0}$ and $Y^{0 A_1^0}$ for some $a_0 \neq 0$. Because A_1^0 may often reflect a natural level of A_1 (as in the absence of treatment) which differs between subjects, natural direct effects may have a more appealing interpretation than controlled direct effects. They moreover correspond with a measure of natural indirect effect in terms of contrasts between $Y^{a_0 A_1^{a_0}}$ and $Y^{a_0 A_1^0}$ for some $a_0 \neq 0$. SMMs for natural direct effects have been considered van der Laan and Petersen (2008) and Tchetgen Tchetgen and Shpitser (2011). Such models are defined by

$$(40) \quad \begin{aligned} E(Y^{a_0 A_1^0} - Y^{0 A_1^0} \mid A_0 = a_0, L_0 = l_0) \\ = m(a_0, l_0; \psi^*), \end{aligned}$$

for each a_0, l_0 , where $m(a_0, L_0; \psi)$ is a known function, smooth in ψ , which satisfies $m(0, L_0; \psi) = 0$. Extensions to sequential treatments or mediators have so far not been developed in view of difficulties of identification in such settings.

9. CONCLUDING REMARKS

Structural nested models were designed in part to deal with confounding by variables affected by treatment. These models maintain close resemblance to ordinary regression models by parameterizing conditional treatment effects. However, in contrast to these, they avoid conditioning on post-treatment variables by modeling the outcome at each time conditional on the treatment and covariate history up to that time; they do this after having removed the effects of later treatments so as to disentangle the unique contributions of each treatment at each time. The associated method of G-estimation has close resemblance to ordinary regression methods because it realizes control for measured confounders through conditioning. In spite of these strong connections with popular estimation methods, SNMs and G-estimation have not become quite as popular as MSMs and the associated IPW methods (Robins, Hernan and Brumback, 2000).

The lack of popularity of G-estimation is largely related to the fact that it cannot usually be performed via off-the-shelf software; however, note that SAS and Stata macros for SNFTMs and SNCFTMs are available at <http://www.hsph.harvard.edu/causal/software/>. This lack of popularity is additionally related to difficulties in solving the estimating equations in the analysis of censored survival times using SNFTMs. These difficulties can now be overcome by using the newer class of SNCFTMs instead (Piciotto et al., 2012; Martinussen et al., 2011).

In spite of these limitations, SNMs and G-estimation allow for greater flexibility than MSMs and typically yield better performing estimators (see Section 4.1). This is especially so when handling continuous exposures or when handling a binary exposure that is strongly correlated with subject characteristics (e.g., when the treated and untreated are very different in terms of subject characteristics). In the latter case, IPW estimators will typically have a poor performance, reflecting the lack of information about the treatment effect in strata where most/all subjects are treated or untreated. In contrast, because SNMs parameterize treatment effects conditionally on covariates, nonsaturated models allow for borrowing of information, so that G-estimators can pool the treatment effects across strata, as in expression (23), downweighing those strata where information on treatment effect is lacking. SNMs can also incorporate effect modification by time-varying covariates. As such, a saturated SNM encodes all possible causal contrasts on the considered scale, in contrast to MSMs which average the effects across (time-varying) covariates, thereby diluting the effects when effect heterogeneity exists on the considered scale. SNMs can moreover make use of instrumental variables.

G-estimation is not to be confused with G-computation (Robins, 1986), which involves standardizing the predictions from an outcome model corresponding to the considered treatment regime, relative to the confounder distribution in the population. Up to recently, also this approach has received little attention in practice because it is computationally intensive and because correct specification of models for the distribution of the (possibly high-dimensional) confounders can be a thorny issue in practice. These concerns, which also relate to likelihood-based inference under SNDMs (Robins, Rotnitzky and Scharfstein, 2000), can be somewhat mitigated by summarising the confounders at each time by a longitudinal propensity score defined as the probability of treatment at that time, given the history of confounders at that time (Achy-Brou, Frangakis and Griswold, 2010). However, this may demand correct specification of propensity score models in addition to a model for the outcome at each time. G-computation moreover does not enable a transparent parameterization of the effect of a particular treatment regime on the outcome and may thereby imply a null paradox (Robins and Wasserman, 1997) according to which tests of the null hypothesis of no effect may be guaranteed to reject in large samples (Robins, 1997). However, recent empirical appli-

cations have turned out to be rather successful (Cain et al., 2011).

We have attempted to make the literature on structural nested models and G-estimation more accessible, while also giving pointers to the related literatures on effect modification and mediation. Variants of SNMs have also been developed to help identify optimal sequences of treatments when treatments may be assigned dynamically as a function of previous treatment and covariate history. In such settings, it is more natural to model the effect of a blip of treatment at m on a particular utility function Y , such as the outcome at the end-of-study time, if all subsequent treatments are optimal; that is, $a_k = a_k^{\text{opt}}(\bar{l}_k, \bar{a}_{k-1})$ for $k > m$. This can be done by parameterizing the so-called regrets: contrasts of $E(Y^{\bar{a}_m, \bar{a}_{m+1}^{\text{opt}}} | \bar{L}_m, \bar{A}_m = \bar{a}_m)$ and $E(Y^{\bar{a}_{m-1}, \bar{a}_m^{\text{opt}}} | \bar{L}_m, \bar{A}_m = \bar{a}_m)$ (Murphy, 2003). Alternatively, since the optimal treatment is unknown, it may be easier to parameterize the effect of a blip of treatment at m relative to no treatment when all future treatments are optimal. This amounts to contrasting $E(Y^{\bar{a}_m, \bar{a}_{m+1}^{\text{opt}}} | \bar{L}_m, \bar{A}_m = \bar{a}_m)$ and $E(Y^{\bar{a}_{m-1}, 0, \bar{a}_{m+1}^{\text{opt}}} | \bar{L}_m, \bar{A}_m = \bar{a}_m)$ (Robins, 2004). We refer the reader to other papers in this issue for detailed accounts of such models. We conclude by expressing our hope that efforts will be continued to develop computational algorithms and corresponding software programs for SNMs, so as to make these methods accessible to a wider audience.

ACKNOWLEDGMENTS

The authors are grateful to the editors and reviewers for very detailed feedback which substantially improved an earlier version of this manuscript. The first author acknowledges support from the Flemish Research Council (FWO) research Grant G.0111.12 and IAP research network Grant no. P07/05 from the Belgian government (Belgian Science Policy). The second author acknowledges support from the US NIH (Grants # R01-DK090385, RC4-MH092722 and R01-MH078016).

REFERENCES

- ACHY-BROU, A. C., FRANGAKIS, C. E. and GRISWOLD, M. (2010). Estimating treatment effects of longitudinal designs using regression models on propensity scores. *Biometrics* **66** 824–833. [MR2758218](#)
- ALMIRALL, D., TEN HAVE, T. and MURPHY, S. A. (2010). Structural nested mean models for assessing time-varying effect moderation. *Biometrics* **66** 131–139. [MR2756699](#)

- BOWDEN, J. and VANSTEELENDT, S. (2011). Mendelian randomization analysis of case-control data using structural mean models. *Stat. Med.* **30** 678–694. [MR2767465](#)
- CAIN, L. E., LOGAN, R., ROBINS, J. M., STERNE, J. A. C., SABIN, C., BANSI, L., JUSTICE, A., GOULET, J., VAN SIGHEM, A., DE WOLF, F., BUCHER, H. C., VON WYL, V., ESTEVE, A., CASABONA, J., DEL AMO, J., MORENO, S., SENG, R., MEYER, L., PEREZ-HOYOS, S., MUGA, R., LODI, S., LANOY, E., COSTAGLIOLA, D. and HERNAN, M. A. (2011). When to initiate combined antiretroviral therapy to reduce mortality and AIDS-defining illness in HIV-infected persons in developed countries: An observational study. *Ann. Intern. Med.* **154** 509–W173.
- CHAMBERLAIN, G. (1987). Asymptotic efficiency in estimation with conditional moment restrictions. *J. Econometrics* **34** 305–334. [MR0888070](#)
- FRANGAKIS, C. E. and RUBIN, D. B. (2002). Principal stratification in causal inference. *Biometrics* **58** 21–29. [MR1891039](#)
- GOETGHELUK, S., VANSTEELENDT, S. and GOETGHEBEUR, E. (2008). Estimation of controlled direct effects. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **70** 1049–1066. [MR2530329](#)
- GREENLAND, S. and ROBINS, J. M. (1986). Identifiability, exchangeability, and epidemiological confounding. *Int. J. Epidemiol.* **15** 412–418.
- GREENLAND, S., ROBINS, J. and PEARL, J. (1999). Confounding and collapsibility in causal inference. *Statist. Sci.* **14** 29–46.
- HENDERSON, R., ANSELL, P. and ALSHIBANI, D. (2010). Regret-regression for optimal dynamic treatment regimes. *Biometrics* **66** 1192–1201. [MR2758507](#)
- HERNÁN, M. A. (2010). The hazards of hazard ratios. *Epidemiology* **21** 13–15.
- HERNÁN, M. A. and ROBINS, J. M. (2006). Instruments for causal inference: An epidemiologist's dream? *Epidemiology* **17** 360–372.
- JOFFE, M. M., SMALL, D. and HSU, C.-Y. (2007). Defining and estimating intervention effects for groups that will develop an auxiliary outcome. *Statist. Sci.* **22** 74–97. [MR2408662](#)
- JOFFE, M. M., YANG, W. P. and FELDMAN, H. I. (2010). Selective ignorability assumptions in causal inference. *Int. J. Biostat.* **6** Art. 11, 25. [MR2602554](#)
- JOFFE, M. M., YANG, W. P. and FELDMAN, H. (2012). G-estimation and artificial censoring: Problems, challenges, and applications. *Biometrics* **68** 275–286. [MR2909884](#)
- MARK, S. D. and ROBINS, J. M. (1993). A method for the analysis of randomized trials with compliance information: An application to the multiple risk factor intervention trial. *Contr. Clin. Trials* **14** 79–97.
- MARTINUSSEN, T., VANSTEELENDT, S., GERSTER, M. and VON BORNE-MANN HJELMBORG, J. (2011). Estimation of direct effects for survival data by using the Aalen additive hazards model. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **73** 773–788. [MR2867458](#)
- MURPHY, S. A. (2003). Optimal dynamic treatment regimes. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **65** 331–366. [MR1983752](#)
- NEWBY, W. K. (1990). Semiparametric efficiency bounds. *J. Appl. Econometrics* **5** 99–135.
- OKUI, R., SMALL, D. S., TAN, Z. and ROBINS, J. M. (2012). Doubly robust instrumental variable regression. *Statist. Sinica* **22** 173–205. [MR2933172](#)
- PEARL, J. (1995). Causal diagrams for empirical research. *Biometrika* **82** 669–710. [MR1380809](#)
- PEARL, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence* 411–420. Morgan Kaufmann, San Francisco, CA.
- PICCIOTTO, S., HERNÁN, M. A., PAGE, J. H., YOUNG, J. G. and ROBINS, J. M. (2012). Structural nested cumulative failure time models to estimate the effects of interventions. *J. Amer. Statist. Assoc.* **107** 886–900. [MR3010878](#)
- ROBINS, J. (1986). A new approach to causal inference in mortality studies with a sustained exposure period—Application to control of the healthy worker survivor effect. Mathematical models in medicine: Diseases and epidemics. Part 2. *Math. Modelling* **7** 1393–1512. [MR0877758](#)
- ROBINS, J. M. (1989). The analysis of randomized and non-randomized AIDS treatment trials using a new approach to causal inference in longitudinal studies. In *Health Service Research Methodology: A Focus on AIDS* (L. Sechrest, H. Freeman and A. Mulley, eds.) 113–159. U.S. Public Health Service, National Center for Health Services Research, Washington, DC.
- ROBINS, J. (1992). Estimation of the time-dependent accelerated failure time model in the presence of confounding factors. *Biometrika* **79** 321–334. [MR1185134](#)
- ROBINS, J. M. (1994). Correcting for non-compliance in randomized trials using structural nested mean models. *Comm. Statist. Theory Methods* **23** 2379–2412. [MR1293185](#)
- ROBINS, J. M. (1997). Causal inference from complex longitudinal data. In *Latent Variable Modeling and Applications to Causality* (Los Angeles, CA, 1994). *Lecture Notes in Statist.* **120** 69–117. Springer, New York. [MR1601279](#)
- ROBINS, J. M. (1999). Testing and estimation of direct effects by reparameterizing directed acyclic graphs with structural nested models. In *Computation, Causation, and Discovery* (C. Glymour and G. Cooper, eds.) 349–405. AAAI Press, Menlo Park, CA. [MR1696459](#)
- ROBINS, J. M. (2000). Marginal structural models versus structural nested models as tools for causal inference. In *Statistical Models in Epidemiology, the Environment, and Clinical Trials* (Minneapolis, MN, 1997) (M. Halloran and D. Berry, eds.) *IMA Vol. Math. Appl.* **116** 95–133. Springer, New York. [MR1731682](#)
- ROBINS, J. M. (2004). Optimal structural nested models for optimal sequential decisions. In *Proceedings of the Second Seattle Symposium in Biostatistics. Lecture Notes in Statist.* **179** 189–326. Springer, New York. [MR2129402](#)
- ROBINS, J. M. and GREENLAND, S. (1992). Identifiability and exchangeability for direct and indirect effects. *Epidemiology* **3** 143–155.
- ROBINS, J. M. and GREENLAND, S. (1994). Adjusting for differential rates of prophylaxis therapy for PCP in high- versus low-dose AZT treatment arms in an AIDS randomized trial. *J. Amer. Statist. Assoc.* **89** 737–749.
- ROBINS, J. M. and HERNÁN, M. A. (2009). Estimation of the causal effects of time-varying exposures. In *Longitudinal Data Analysis* (G. Fitzmaurice, M. Davidian, G. Verbeke and G. Molenberghs, eds.) 553–599. CRC Press, Boca Raton, FL. [MR1500133](#)
- ROBINS, J. M., HERNÁN, M. A. and BRUMBACK, B. (2000). Marginal structural models and causal inference in epidemiology. *Epidemiology* **11** 550–560.

- ROBINS, J. M., MARK, S. D. and NEWBY, W. K. (1992). Estimating exposure effects by modelling the expectation of exposure conditional on confounders. *Biometrics* **48** 479–495. [MR1173493](#)
- ROBINS, J. M. and RITOV, Y. (1997). Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semiparametric models. *Stat. Med.* **16** 285–319.
- ROBINS, J. M. and ROTNITZKY, A. (2001). Comment on “Inference for semiparametric models: Some questions and an answer,” by P. J. Bickel and J. Kwon. *Statist. Sinica* **11** 920–936.
- ROBINS, J. and ROTNITZKY, A. (2004). Estimation of treatment effects in randomised trials with non-compliance and a dichotomous outcome using structural mean models. *Biometrika* **91** 763–783. [MR2126032](#)
- ROBINS, J. M., ROTNITZKY, A. and SCHARFSTEIN, D. O. (2000). Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In *Statistical Models in Epidemiology, the Environment, and Clinical Trials* (Minneapolis, MN, 1997) (M. Halloran and D. Berry, eds.), *IMA Vol. Math. Appl.* **116** 1–94. Springer, New York. [MR1731681](#)
- ROBINS, J. M., ROTNITZKY, A. and ZHAO, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed. *J. Amer. Statist. Assoc.* **89** 846–866. [MR1294730](#)
- ROBINS, J. M. and TSIATIS, A. A. (1991). Correcting for non-compliance in randomized trials using rank preserving structural failure time models. *Comm. Statist. Theory Methods* **20** 2609–2631. [MR1144866](#)
- ROBINS, J. M. and WASSERMAN, L. (1997). Estimation of Effects of Sequential Treatments by Reparameterizing Directed Acyclic Graphs. In *Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence* (D. Geiger and P. Shenoy, eds.) 409–420. Morgan Kaufmann, San Francisco, CA.
- ROBINS, J. M., BLEVINS, D., RITTER, G. and WULFSOHN, M. (1992). G-estimation of the effect of prophylaxis therapy for pneumocystic carinii pneumonia on the survival of AIDS patients. *Epidemiology* **3** 319–336.
- ROSENBAUM, P. R. (1984). The consequences of adjustment for a concomitant variable that has been affected by the treatment. *J. Roy. Statist. Soc. Ser. A* **147** 656–666.
- ROSENBAUM, P. R. and RUBIN, D. B. (1984). Reducing bias in observational studies using subclassification on the propensity score. *J. Amer. Statist. Assoc.* **79** 516–524.
- RUBIN, D. B. (1978). Bayesian inference for causal effects: The role of randomization. *Ann. Statist.* **6** 34–58. [MR0472152](#)
- STEPHENS, A., KEELE, L. and JOFFE, M. (2013). Estimating post-treatment effect modification with generalized structural mean models. Submitted.
- TCHETGEN TCHETGEN, E. J. (2012). Multiple-robust estimation of an odds ratio interaction. Harvard Univ. Biostatistics working paper series. Working Paper 142. Available at <http://biostats.bepress.com/harvardbiostat/paper142>.
- TCHETGEN TCHETGEN, E. J. and ROBINS, J. (2010). The semiparametric case-only estimator. *Biometrics* **66** 1138–1144. [MR2758501](#)
- TCHETGEN TCHETGEN, E. J., ROBINS, J. M. and ROTNITZKY, A. (2010). On doubly robust estimation in a semiparametric odds ratio model. *Biometrika* **97** 171–180. [MR2594425](#)
- TCHETGEN TCHETGEN, E. J. and ROTNITZKY, A. (2011). Double-robust estimation of an exposure-outcome odds ratio adjusting for confounding in cohort and case-control studies. *Stat. Med.* **30** 335–347. [MR2758866](#)
- TCHETGEN TCHETGEN, E. J. and SHPITSER, I. (2014). Estimation of a semiparametric natural direct effect model incorporating baseline covariates. *Biometrika* **101** 849–864.
- TEN HAVE, T. R., JOFFE, M. M., LYNCH, K. G., BROWN, G. K., MAISTO, S. A. and BECK, A. T. (2007). Causal mediation analyses with rank preserving models. *Biometrics* **63** 926–934. [MR2395813](#)
- VANSTEELENDT, S. (2010). Estimation of controlled direct effects on a dichotomous outcome using logistic structural direct effect models. *Biometrika* **97** 921–934. [MR2746161](#)
- VANSTEELENDT, S., BEKAERT, M. and CLAESKENS, G. (2012). On model selection and model misspecification in causal inference. *Stat. Methods Med. Res.* **21** 7–30. [MR2867536](#)
- VANSTEELENDT, S. and DANIEL, R. M. (2014). On regression adjustment for the propensity score. *Stat. Med.* **33** 4053–4072.
- VANSTEELENDT, S. and GOETGHEBEUR, E. (2003). Causal inference with generalized structural mean models. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **65** 817–835. [MR2017872](#)
- VANSTEELENDT, S., VANDERWEELE, T., TCHETGEN, E. J. and ROBINS, J. M. (2008a). Semiparametric inference for statistical interactions. *J. Amer. Statist. Assoc.* **103** 1693–1704.
- VANSTEELENDT, S., DEMEO, D. L., LASKY-SU, J. et al. (2008b). Testing and estimating gene-environment interactions in family-based association studies. *Biometrics* **64** 458–467. [MR2432416](#)
- VANSTEELENDT, S., BOWDEN, J., BABANEZHAD, M. and GOETGHEBEUR, E. (2011). On instrumental variables estimation of causal odds ratios. *Statist. Sci.* **26** 403–422. [MR2917963](#)
- VAN DER LAAN, M. J., HUBBARD, A. and JEWELL, N. P. (2007). Estimation of treatment effects in randomized trials with non-compliance and a dichotomous outcome. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **69** 463–482. [MR2323763](#)
- VAN DER LAAN, M. J. and PETERSEN, M. L. (2008). Direct effect models. *Int. J. Biostat.* **4** 1–27. [MR2456975](#)
- VOCK, D. M., TSIATIS, A. A., DAVIDIAN, M., LABER, E. B., TSUANG, W. M., FINLEN COPELAND, C. A. and PALMER, S. M. (2013). Assessing the causal effect of organ transplantation on the distribution of residual lifetime. *Biometrics* **69** 820–829. [MR3146778](#)
- WEI, L. J. (1992). The accelerated failure time model: A useful alternative to the Cox regression model in survival analysis. *Stat. Med.* **11** 1871–1879.
- ZHANG, M., JOFFE, M. M. and SMALL, D. S. (2011). Causal inference for continuous-time processes when covariates are observed only at discrete times. *Ann. Statist.* **39** 131–173. [MR2797842](#)